

---

# *Linear Models of Classification: Probabilistic Generative and Discriminative Models*

*Prof. Nicholas Zabaras*

*Center for Informatics and Computational Science*

<https://cics.nd.edu/>

*University of Notre Dame  
Notre Dame, Indiana, USA*

*Email: [nzabaras@gmail.com](mailto:nzabaras@gmail.com)*

*URL: <https://www.zabaras.com/>*

*February 19, 2019*

# Contents

---

- Linear models for classification, Approaches to Classification, Discriminant Functions, Least Squares approach to Classification
- Generative Vs Discriminative Classifiers, Fishers linear discriminant, Probabilistic Interpretation
- Probabilistic Generative Models for two Classes, Logistic Sigmoid, Models with K>2, Gaussian Class Conditionals, MLE, Multiclass Case, Discrete Features, Generalization to Exponential Family Class Conditionals
- Online Learning and Stochastic Optimization – Robbins-Monro Algorithm, The Perceptron algorithm, Bayesian perspective
- Probabilistic Discriminative Models, Nonlinear Basis in Classification Models, Logistic Regression and Generalized Linear Models, Cross Entropy Error, Sequential Update, Linearly Separable Data, Regularization, Iterative Reweighted Least Squares, Multiclass Logistic Regression
- Algorithms for Unconstrained Optimization, Steepest Descent Method, Newton's Method, BFGS, Regularization
- Bayesian Logistic Regression, Laplace Approximation, Model comparison and Model Evidence, BIC Criterion, Gaussian Approximation to the Posterior, Predictive Distribution, Probit Approximation, Outlier Detection
- Probit Regression, IRLS for Probit Regression, Canonical Link Functions
  - Following closely Chris Bishop's PRML book, Chapter 4.
  - K. Murphy, Machine Learning: A probabilistic Perspective, Chapter 8

# Linear Models for Classification

---

- The goal in classification is to take a  $D$  –dimensional input  $\mathbf{x}$  and to assign it to one of  $K$  discrete classes  $C_k, k = 1, \dots, K$ .
- **Decision Regions and Decision Boundaries:** The input space is divided into **decision regions** whose boundaries are called **decision boundaries**.
- **Linear Models for Classification:** The decision boundaries are linear functions of  $\mathbf{x}$  defined by  $(D - 1)$  –dimensional hyperplanes within the  $D$  –dimensional input space.
- **Linearly Separable Data Sets:** Data sets whose classes can be separated exactly by linear decision surfaces are said to be **linearly separable**.

# Classification: Target Values

---

- There are many ways of using target values to represent class labels.
- For probabilistic 2 class models, there is a single target variable  $t \in \{0,1\}$ .  $t = 1$  represents class  $C_1$ ,  $t = 0$  class  $C_2$ .
  - We interpret the value of  $t$  as the probability that the class is  $C_1$ .
- For  $K > 2$  classes, we use a **1-of- $K$  coding** in which  $\mathbf{t}$  is a vector of length  $K$  such that if the class is  $C_j$ , then all elements  $t_k$  of  $\mathbf{t}$  are zero except  $t_j = 1$ .
  - E.g., for  $K = 5$ , then a pattern from class 2 has the target
$$\mathbf{t} = (0, 1, 0, 0, 0)^T$$
  - We interpret the value of  $t_k$  as the probability that the class is  $C_k$ .

# Approaches to Classification

- **Case I:** Construct a discriminant function that directly assigns each vector  $x$  to a specific class.
- **Case II:** Model the conditional probability distribution  $p(C_k|x)$  in an inference stage, and subsequently use this to make optimal decisions.
  - **Case IIa:** Model  $p(C_k|x)$  directly by parametric models and then optimize the parameters using a training set.
  - **Case IIb:** Use a generative approach in which we model the class-conditional densities  $p(x|C_k)$ , together with the priors  $p(C_k)$  for the classes, and then compute  $p(C_k|x)$  using Bayes' theorem

$$p(C_k | x) = \frac{p(x | C_k)p(C_k)}{p(x)}$$

---

# **Discriminant Functions**

# Generalized Linear Models of Classification

---

- To predict discrete class labels e.g. posterior probabilities in the range  $(0, 1)$ , we consider *a generalization of the* Bayesian regression model by transforming a linear function of  $w$  with a nonlinear function  $f(\cdot)$ :

$$y(\mathbf{x}) = f(\mathbf{w}^T \mathbf{x} + w_0)$$

- $f(\cdot)$  is known as the **activation function**.
- The decision surfaces correspond to  $y(\mathbf{x}) = \text{constant}$ , so that  $\mathbf{w}^T \mathbf{x} + w_0 = \text{const.}$ , i.e. **the decision surfaces are linear in  $x$ , even if  $f(\cdot)$  is nonlinear.**
- These models are called *generalized linear models*.

▪ McCullagh, P. and J. Nelder (1989). *Generalized linear models*. Chapman and Hall. 2nd edition.

# Generalized Linear Models of Classification

---

- Note that the generalized linear models of classification

$$y(\mathbf{x}) = f(\mathbf{w}^T \mathbf{x} + w_0)$$

are not linear in  $\mathbf{w}$  as is the case in regression models.

- As a result, classifications models are more complex analytically and computationally.
- The classification techniques can take place
  - directly in the input space  $\mathbf{x}$  or
  - in a feature space defined by the basis functions  $\phi(\mathbf{x})$ .

# Discriminant Functions

---

- A discriminant is a function that takes  $x$  and assigns it directly to one of the  $K$  classes  $C_k$ .
- We consider here linear discriminants for which the decision surfaces are hyperplanes.
- For two classes, a linear discriminant function is of the form:

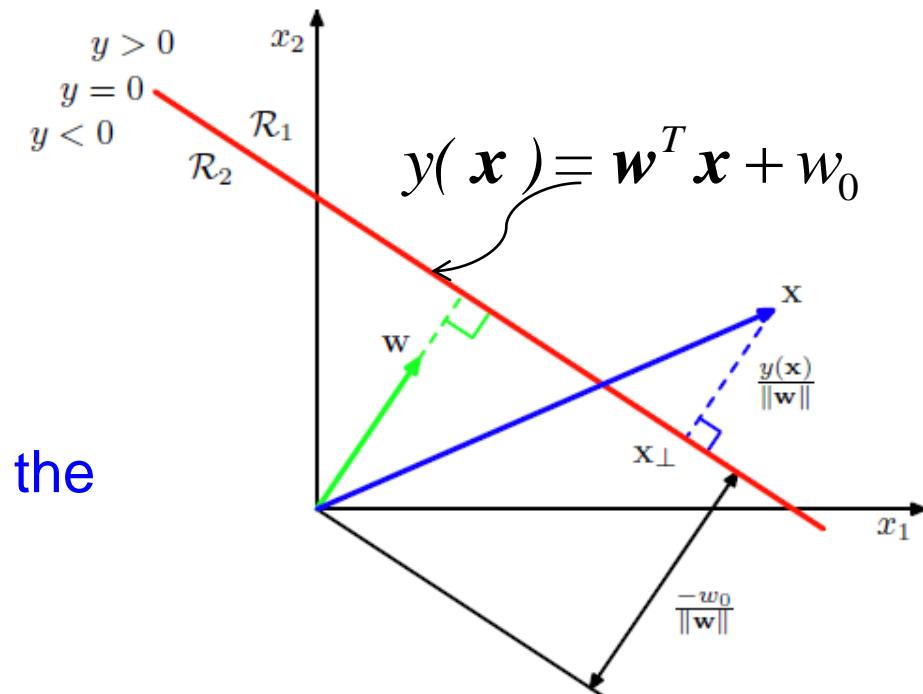
$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

$\mathbf{w}$  is the weight vector,  $w_0$  the bias ( $-w_0$  is called the threshold).

- We consider next the case of two classes.

# Discriminant Functions For Two Classes

- $x$  is assigned to  $C_1$  if  $y(x) \geq 0$  and to  $C_2$  otherwise.
- The decision boundary is  $y(x) = 0$ , a  $(D - 1)$ -dimensional hyperplane within the  $D$ -dimensional input space.
- Consider two points  $x_A$  and  $x_B$  on the decision surface. Because  $y(x_A) = y(x_B) = 0$ , we have  $\mathbf{w}^T(x_A - x_B) = 0$  and hence  $\mathbf{w}$  is orthogonal to every vector on the decision surface.
- $\mathbf{w}$  determines the orientation of the decision surface.

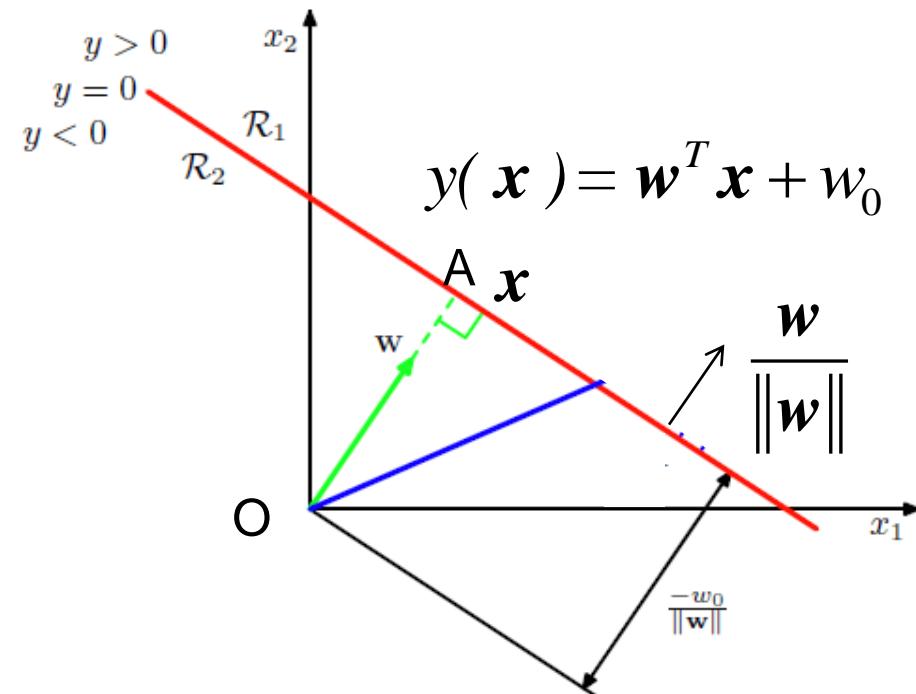


# Discriminant Functions For Two Classes

- Consider the normal distance (OA) from the origin to the decision surface. Let  $x$  the projection of the origin on the decision surface. Since  $y(x) = 0$ ,

$$\overrightarrow{OA} = x = \alpha \frac{\mathbf{w}}{\|\mathbf{w}\|} \Rightarrow \mathbf{w}^T \alpha \frac{\mathbf{w}}{\|\mathbf{w}\|} + w_0 = 0 \Rightarrow (OA) = \alpha = -\frac{w_0}{\|\mathbf{w}\|}$$

- The bias parameter  $w_0$  determines the location of the decision surface from the origin.



# Discriminant Functions For Two Classes

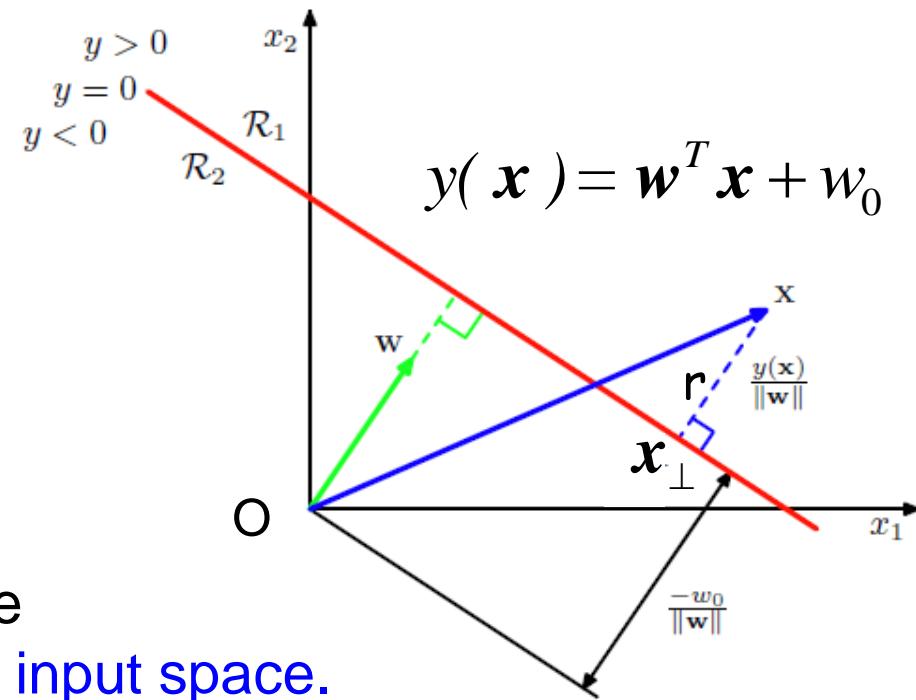
- We note that the value of  $y(\mathbf{x})$  for any  $\mathbf{x}$  gives a signed measure of the distance  $r$  of  $\mathbf{x}$  from the decision surface.

$$\mathbf{x} = \mathbf{x}_\perp + r \frac{\mathbf{w}}{\|\mathbf{w}\|} \Rightarrow \underbrace{\mathbf{w}^T \mathbf{x} + w_0}_{y(\mathbf{x})} = \underbrace{\mathbf{w}^T \mathbf{x}_\perp + w_0}_{=0} + r \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|} = r \|\mathbf{w}\| \Rightarrow r = \frac{y(\mathbf{x})}{\|\mathbf{w}\|}$$

- We can use compact notation and introduce  $x_0 = 1$  and

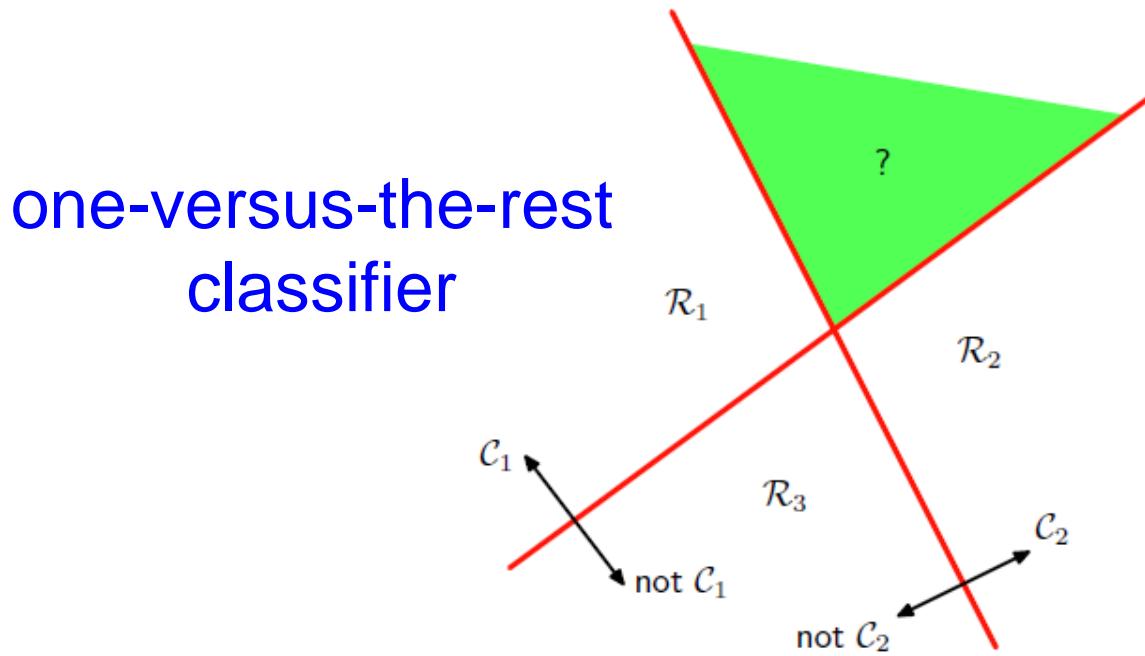
$$\tilde{\mathbf{w}} = (w_0, \mathbf{w}), \tilde{\mathbf{x}} = (x_0, \mathbf{x}) \\ y(\tilde{\mathbf{x}}) = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}$$

- The decision surfaces are now  $D$  –dimensional hyperplanes passing through the origin of the  $D + 1$  –dimensional expanded input space.



# Discriminant Functions For Many Classes

- Building a  $K$  –class discriminant by combining a number of two-class discriminant functions does not work in general.
- Consider the **use of  $(K - 1)$  classifiers** each of which solves a two-class problem of **separating points** in a particular class  $C_k$  from points not in that class.

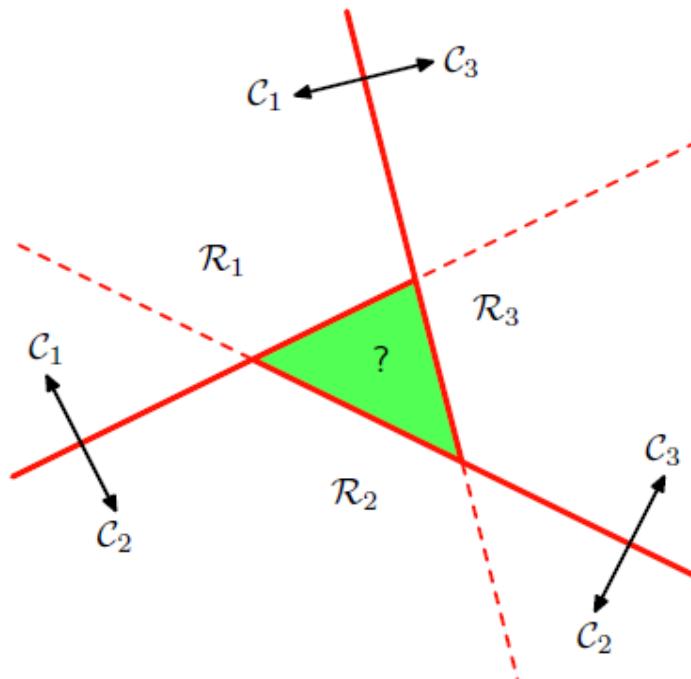


- Duda, R. O., P. E. Hart, and D. G. Stork (2001). *Pattern Classification*. Wiley Interscience. 2nd edition.

# Discriminant Functions For Many Classes

- One can introduce  $K(K - 1)/2$  binary discriminant functions, one for every possible pair of classes.
- Each point is classified using a majority vote amongst the discriminant functions.
- This also leads to ambiguous regions.

one-versus-one  
classifier



# **Discriminant Functions For Many Classes**

---

- To resolve the problems of the earlier classifiers, we consider a single  $K$  –class discriminant comprising  $K$  linear functions of the form

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}$$

- We assign  $\mathbf{x}$  to class  $C_k$  if  $y_k(\mathbf{x}) > y_j(\mathbf{x})$  for all  $j \neq k$ .
- The decision boundary between class  $C_k$  and class  $C_j$  is therefore given by  $y_k(\mathbf{x}) = y_j(\mathbf{x})$ . It is a  $(D - 1)$  –dimensional hyperplane

$$(\mathbf{w}_k - \mathbf{w}_j)^T \mathbf{x} + w_{k0} - w_{j0} = 0$$

- It has the same form as the decision boundary for the two-class case.

# Discriminant Functions For Many Classes

- The decision regions of this discriminant are singly connected and convex.
- Consider  $\mathbf{x}_A$  and  $\mathbf{x}_B$  which lie inside decision region  $\mathcal{R}_k$ . Any  $\mathbf{x}$  on the line connecting  $\mathbf{x}_A$  and  $\mathbf{x}_B$  can be expressed as

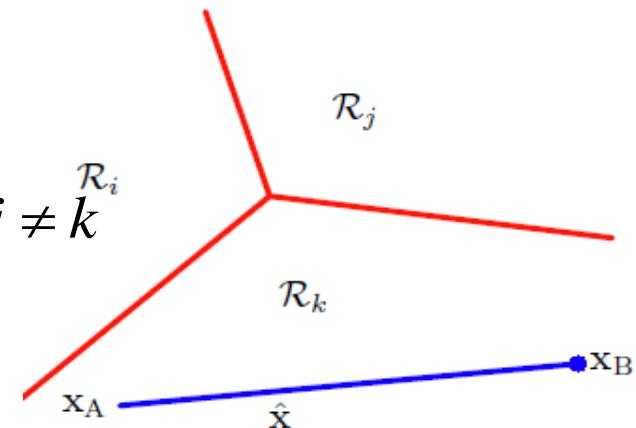
$$\hat{\mathbf{x}} = \lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B, \quad 0 \leq \lambda \leq 1$$

- From the linearity of the discriminant functions we have:

$$y_k(\hat{\mathbf{x}}) = \lambda y_k(\mathbf{x}_A) + (1 - \lambda) y_k(\mathbf{x}_B)$$

- Using  $y_k(\mathbf{x}_A) > y_j(\mathbf{x}_A)$ ,  $y_k(\mathbf{x}_B) > y_j(\mathbf{x}_B)$ ,  $j \neq k$

we see:  $y_k(\hat{\mathbf{x}}) > y_j(\hat{\mathbf{x}})$ ,  $j \neq k$  and so  $\hat{\mathbf{x}} \in \mathcal{R}_k$ .



- For 2 classes, one can use 2 discriminants or a single one as discussed earlier.

---

# *Least Squares for Classification*

# Least Squares Formalism to Classification

- Consider a general classification problem with  $K$  classes, with a 1-of- $K$  binary coding scheme for the target  $t$ .
- One justification for *the least squares approach is that it approximates the conditional expectation  $\mathbb{E}[t|x]$ .*
  - For the binary coding scheme,  $\mathbb{E}[t|x]$  is given by the vector of posterior class probabilities.
  - The approximation of these posterior class probabilities is poor and can have values outside the range  $(0, 1)$ .

# Least Squares Formalism to Classification

- Each class is defined by a linear model:

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}, k = 1, \dots, K$$

- We put all these equations together in a vector form:

$$\mathbf{y}(\mathbf{x}) = \widetilde{\mathbf{W}}^T \widetilde{\mathbf{x}}$$

- The  $k^{\text{th}}$  column of  $\widetilde{\mathbf{W}}$  ( $D + 1 \times K$ ) comprises the  $D + 1$ -dim vector  $\widetilde{\mathbf{w}}_k = (w_{k0}, \mathbf{w}_k^T)^T$   
and  $\widetilde{\mathbf{x}} = (1, \mathbf{x}^T)^T$ .

$$\widetilde{\mathbf{W}} = \left( \begin{array}{c} \uparrow \\ \widetilde{\mathbf{w}}_k \\ \downarrow \end{array} \right)$$

- A new input  $\mathbf{x}$  is then assigned to the class for which the output  $y_k(\mathbf{x}) = \widetilde{\mathbf{w}}_k^T \widetilde{\mathbf{x}}$  is the largest.

# Least Squares Formalism to Classification

- Consider a training data set:  $\{\mathbf{x}_n, \mathbf{t}_n\}, n = 1, \dots, N$ .

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}, k = 1, \dots, K \quad \text{or} \quad \mathbf{y}(\mathbf{x}) = \widetilde{\mathbf{W}}^T \widetilde{\mathbf{x}}$$

- We define matrices  $\widetilde{\mathbf{W}}, \mathbf{T}, \widetilde{\mathbf{X}}$  as follows:

$$\widetilde{\mathbf{X}} = \left( \begin{array}{c} \leftarrow \widetilde{\mathbf{x}}_n^T \rightarrow \\ \end{array} \right) \text{n}^{\text{th}} \text{ row} \quad \widetilde{\mathbf{W}} = \left( \begin{array}{c} \uparrow \\ \widetilde{\mathbf{w}}_k \\ \downarrow \end{array} \right)_{(D+1) \times K} \quad \mathbf{T} = \left( \begin{array}{c} \leftarrow \mathbf{t}_n^T \rightarrow \\ \end{array} \right) \text{n}^{\text{th}} \text{ row} \\ N \times (D+1) \qquad \qquad \qquad N \times K$$

- We can now write the **sum-of-squares error** and minimize:

$$E_D(\widetilde{\mathbf{W}}) = \frac{1}{2} \text{Tr} \left\{ (\widetilde{\mathbf{X}} \widetilde{\mathbf{W}} - \mathbf{T})^T (\widetilde{\mathbf{X}} \widetilde{\mathbf{W}} - \mathbf{T}) \right\} \Rightarrow \widetilde{\mathbf{W}} = (\widetilde{\mathbf{X}}^T \widetilde{\mathbf{X}})^{-1} \widetilde{\mathbf{X}}^T \mathbf{T} = \widetilde{\mathbf{X}}^\dagger \mathbf{T}$$

$$\text{Note : } \text{tr}(\mathbf{A}^T \mathbf{A}) = \sum_{i,j} A_{ij}^2$$

# Least Squares Formalism to Classification

- For a new input  $x$ , we obtain the following discriminant function:

$$y(x) = \widetilde{W}^T \tilde{x} = T^T (\tilde{X}^\dagger)^T \tilde{x}$$

- Note that if every target vector in the training set satisfies some linear constraint  $a^T t_n + b = 0$  for some  $a, b$ , then the model prediction for any  $x$  satisfies ([see proof here](#))

$$a^T y(x) + b = 0$$

- If we use a 1-of-K coding scheme for  $K$  classes, the elements of  $y(x)$  will sum to 1 for any  $x$ .
- However, this constraint is not sufficient to allow the model outputs to be interpreted as probabilities since they are not constrained to lie within  $(0, 1)$ .

# Appendix: Least Squares For Classification

- We can re-write the earlier equations identifying the bias:

$$E_D(\widetilde{\mathbf{W}}) = \frac{1}{2} \text{Tr} \left\{ (\widetilde{\mathbf{X}}\widetilde{\mathbf{W}} - \mathbf{T})^T (\widetilde{\mathbf{X}}\widetilde{\mathbf{W}} - \mathbf{T}) \right\} = \frac{1}{2} \text{Tr} \{ (\mathbf{X}\mathbf{W} + \mathbf{1}\mathbf{w}_0^T - \mathbf{T})^T (\mathbf{X}\mathbf{W} + \mathbf{1}\mathbf{w}_0^T - \mathbf{T}) \}$$

- $\mathbf{w}_0$  is the column vector of the bias weights (top row transposed of  $\widetilde{\mathbf{W}}$ ). Differentiation wrt  $\mathbf{w}_0 \rightarrow$  a familiar result:\*

$$2N\mathbf{w}_0 + 2(\mathbf{X}\mathbf{W} - \mathbf{T})^T \mathbf{1} \Rightarrow \mathbf{w}_0 = \bar{\mathbf{t}} - \mathbf{W}^T \bar{\mathbf{x}}, \quad \bar{\mathbf{t}} = \frac{1}{N} \mathbf{T}^T \mathbf{1}, \quad \bar{\mathbf{x}} = \frac{1}{N} \mathbf{X}^T \mathbf{1}$$

- The error function now takes the form:

$$E_D(\mathbf{W}) = \frac{1}{2} \text{Tr} \{ (\mathbf{X}\mathbf{W} + \bar{\mathbf{T}} - \bar{\mathbf{X}}\mathbf{W} - \mathbf{T})^T (\mathbf{X}\mathbf{W} + \bar{\mathbf{T}} - \bar{\mathbf{X}}\mathbf{W} - \mathbf{T}) \}, \quad \bar{\mathbf{T}} = \mathbf{1}\bar{\mathbf{t}}^T, \quad \bar{\mathbf{X}} = \mathbf{1}\bar{\mathbf{x}}^T$$

- Setting the derivative wrt  $\mathbf{W}$  equal to zero gives:

$$\mathbf{W} = (\widehat{\mathbf{X}}^T \widehat{\mathbf{X}})^{-1} \widehat{\mathbf{X}}^T \widehat{\mathbf{T}} = \widehat{\mathbf{X}}^\dagger \widehat{\mathbf{T}}, \quad \widehat{\mathbf{X}} = \mathbf{X} - \bar{\mathbf{X}}, \quad \widehat{\mathbf{T}} = \mathbf{T} - \bar{\mathbf{T}}$$

\* Use the identities  $\frac{\partial}{\partial \mathbf{a}} \text{Tr}(\mathbf{a}\mathbf{a}^T) = \mathbf{a}$  and  $\frac{\partial}{\partial \mathbf{a}} \text{Tr}(\mathbf{a}\mathbf{b}^T) = \mathbf{b}$ . Thus e.g.  $\frac{\partial}{\partial \mathbf{w}_0} \text{Tr} \left( (\mathbf{1}\mathbf{w}_0^T)^T (\mathbf{1}\mathbf{w}_0^T) \right) = N \frac{\partial}{\partial \mathbf{w}_0} \text{Tr}(\mathbf{w}_0 \mathbf{w}_0^T) = N\mathbf{w}_0$ .

# Appendix: Least Squares For Classification

- For a new input  $x^*$ , the prediction is:

$$y(x^*) = \mathbf{W}^T x^* + \mathbf{w}_0 = \mathbf{W}^T x^* + \bar{\mathbf{t}} - \mathbf{W}^T \bar{x} = \bar{\mathbf{t}} - \widehat{\mathbf{T}}^T (\widehat{\mathbf{X}}^\dagger)^T (x^* - \bar{x})$$

- Consider for some  $\mathbf{a}, b$ ,  $\mathbf{a}^T \mathbf{t}_n + b = 0$ . Apply this to  $\bar{\mathbf{t}} = \frac{1}{N} \mathbf{T}^T \mathbf{1}$

$$\mathbf{a}^T \bar{\mathbf{t}} = \frac{1}{N} \mathbf{a}^T \mathbf{T}^T \mathbf{1} = -b \quad \mathbf{T} = \begin{pmatrix} & & \\ \overleftarrow{t_n^T} & & \\ & & \end{pmatrix}_{N \times K}$$

- Applying  $\mathbf{a}^T \mathbf{t}_n + b = 0$  to  $y(x^*) = \bar{\mathbf{t}} - \widehat{\mathbf{T}}^T (\widehat{\mathbf{X}}^\dagger)^T (x^* - \bar{x})$  we obtain:

$$\mathbf{a}^T y(x^*) = \mathbf{a}^T \bar{\mathbf{t}} - \mathbf{a}^T \widehat{\mathbf{T}}^T (\widehat{\mathbf{X}}^\dagger)^T (x^* - \bar{x}) = \mathbf{a}^T \bar{\mathbf{t}} = -b$$

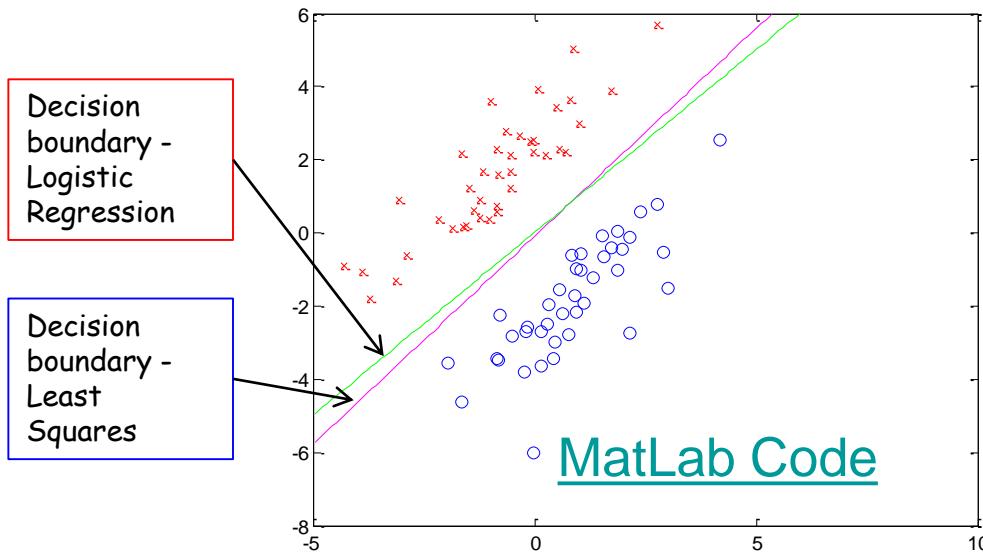
where  $\mathbf{a}^T \widehat{\mathbf{T}}^T = \mathbf{a}^T (\mathbf{T} - \bar{\mathbf{T}})^T = \mathbf{a}^T \mathbf{T}^T - \mathbf{a}^T \bar{\mathbf{T}} \mathbf{1}^T = -b \mathbf{1}^T + b \mathbf{1}^T = \mathbf{0}^T$

- This is the proof of the earlier remark.

# Least Squares Formalism to Classification

- The discriminant function :  $y(x) = \tilde{\mathbf{W}}^T \tilde{\mathbf{x}} = \mathbf{T}^T (\tilde{\mathbf{X}}^\dagger)^T \tilde{\mathbf{x}}$  can be used to make decisions directly without any probabilistic interpretation.
- The least-squares solution in the example below defines two activation functions.
  - The decision boundary (for least squares) shown is taken as:

$$(\mathbf{w}_1 - \mathbf{w}_2)^T \mathbf{x} + w_{10} - w_{20} = 0$$

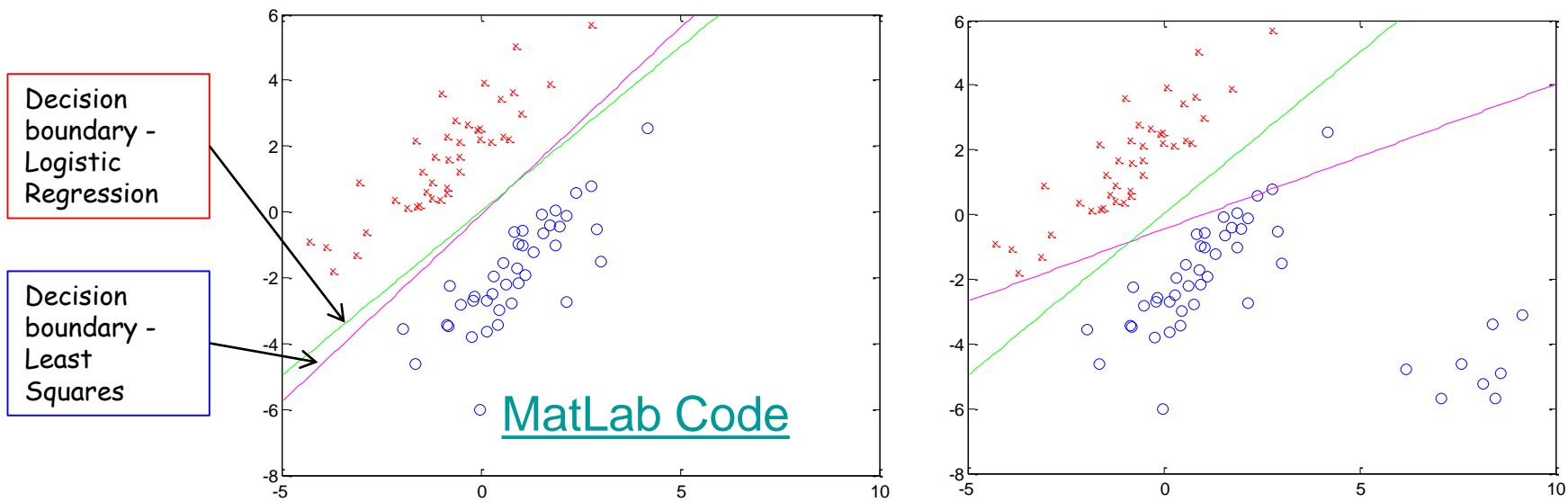


- We will see in a follow up lecture that the decision boundary for the 2 – class logistic regression algorithm in the input space  $\mathbf{x}$  is defined as

$$p(C_1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x}) = 0.5 \Rightarrow \mathbf{w}^T \mathbf{x} = 0$$

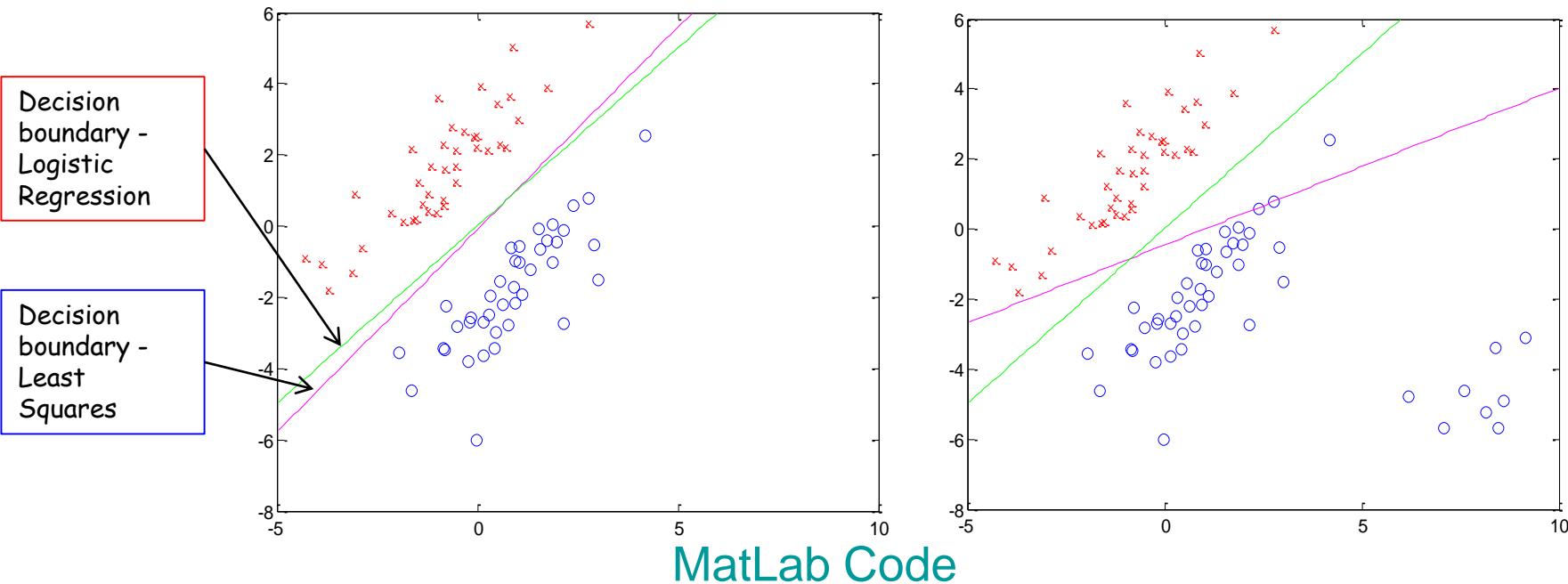
# Least Squares Formalism to Classification

- As is the case for regression problems, the least-squares solution lacks robustness to outliers (see Fig).
  - The sum-of-squares error function penalizes predictions that are ‘too correct’ in that they lie a long way on the correct side of the decision boundary.



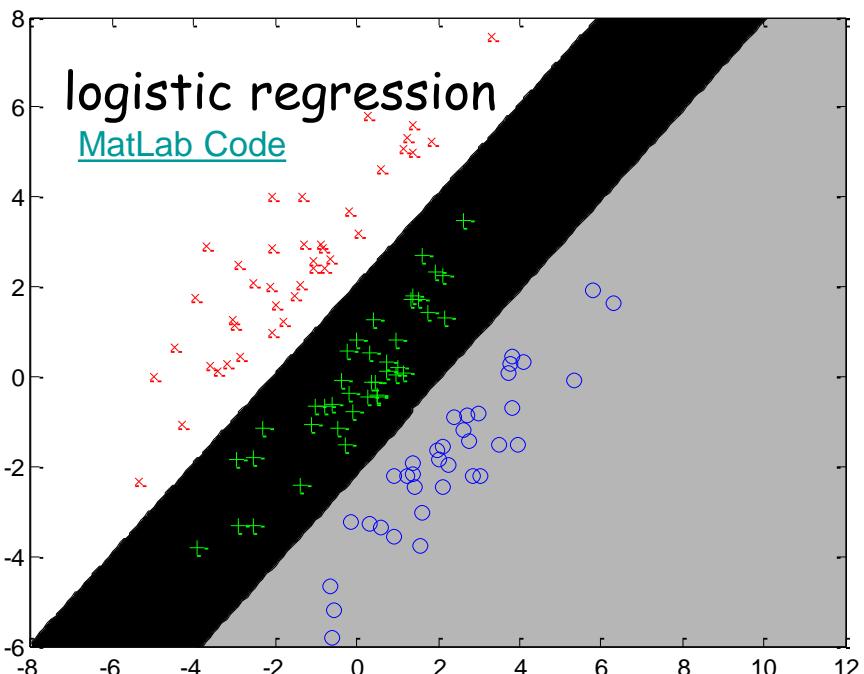
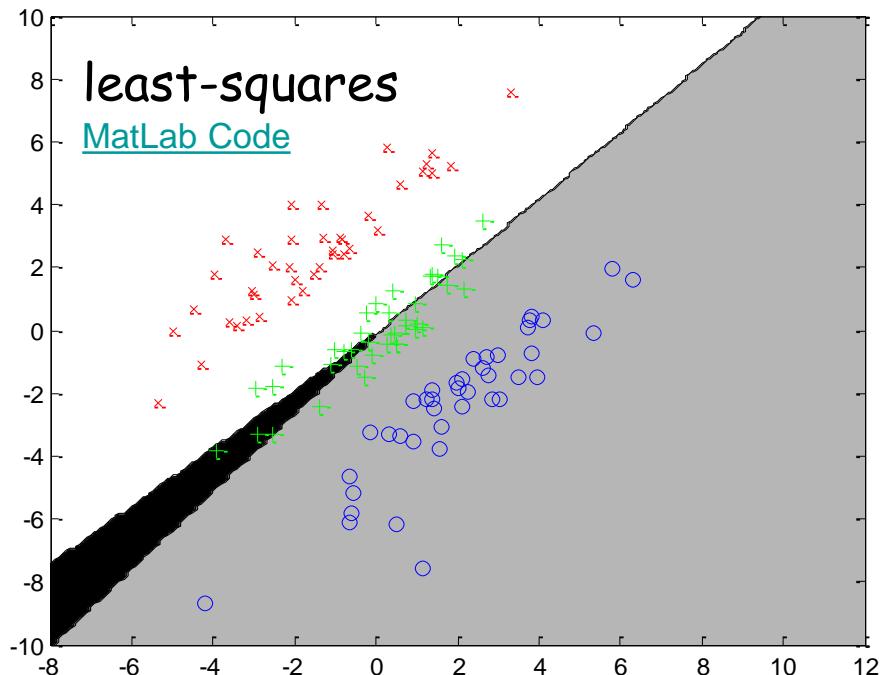
# Least Squares is Sensitive to Outliers

- Additional data points (right) produce a significant change in the location of the decision boundary (compare with the results from logistic regression in the input space  $x$  to be discussed in a coming lecture)
  - Note that these points would be correctly classified by the original decision boundary on the left.



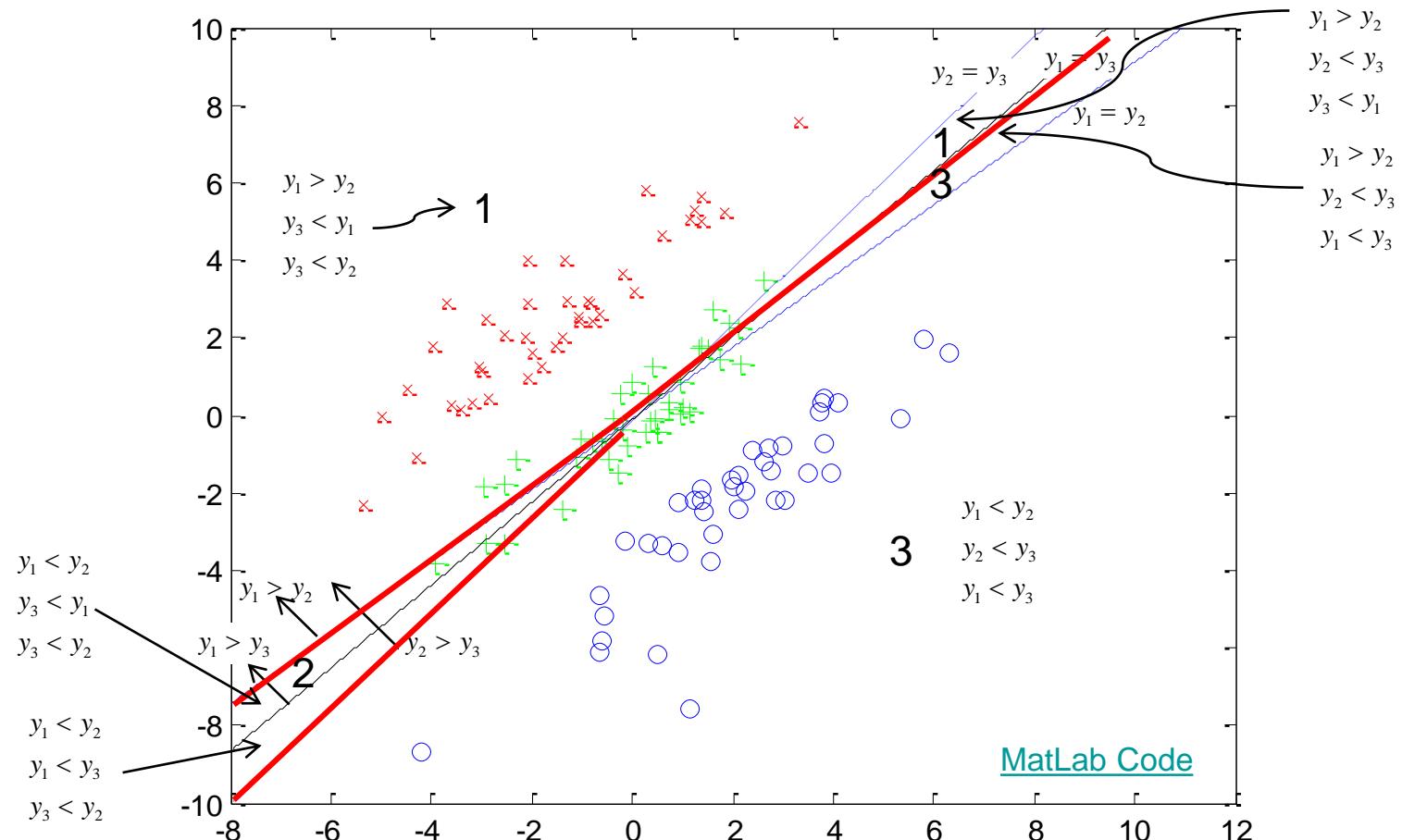
# Least Squares for Multiclass Classification

- Here we consider three classes.
- Three activation functions are defined and for each point, class  $k$  is assigned if  $y_k(x) > y_j(x)$  for all  $j \neq k$ .
- For logistic regression, the regions are assigned based on posterior probabilities  $p(C_m|\phi) = \max(p(C_1|\phi), p(C_2|\phi), p(C_3|\phi))$  with  $p(C_k|\phi) = \exp(a_k)/\sum_j \exp(a_j)$ ,  $a_k = w_k^T \phi$



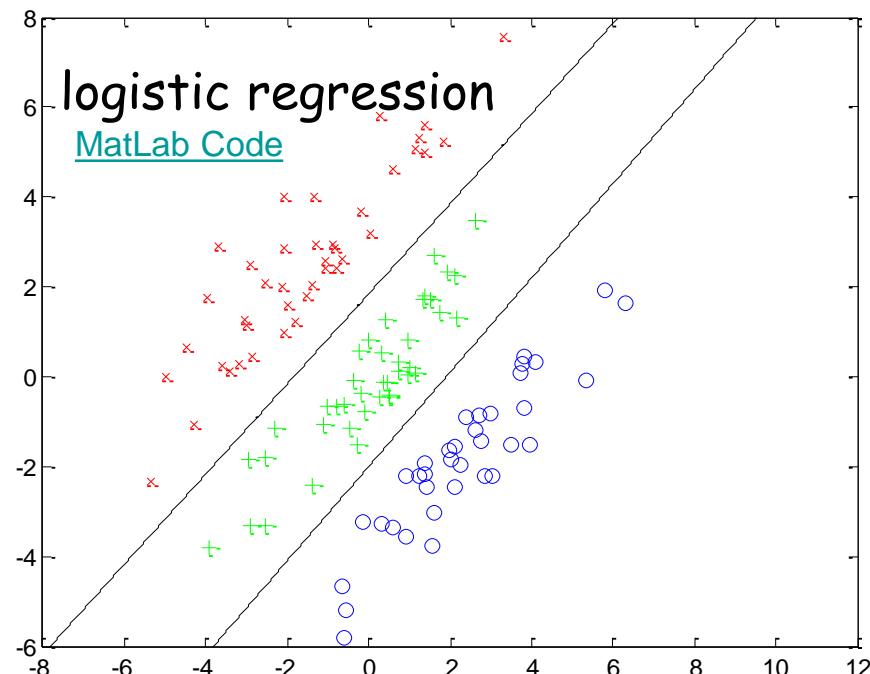
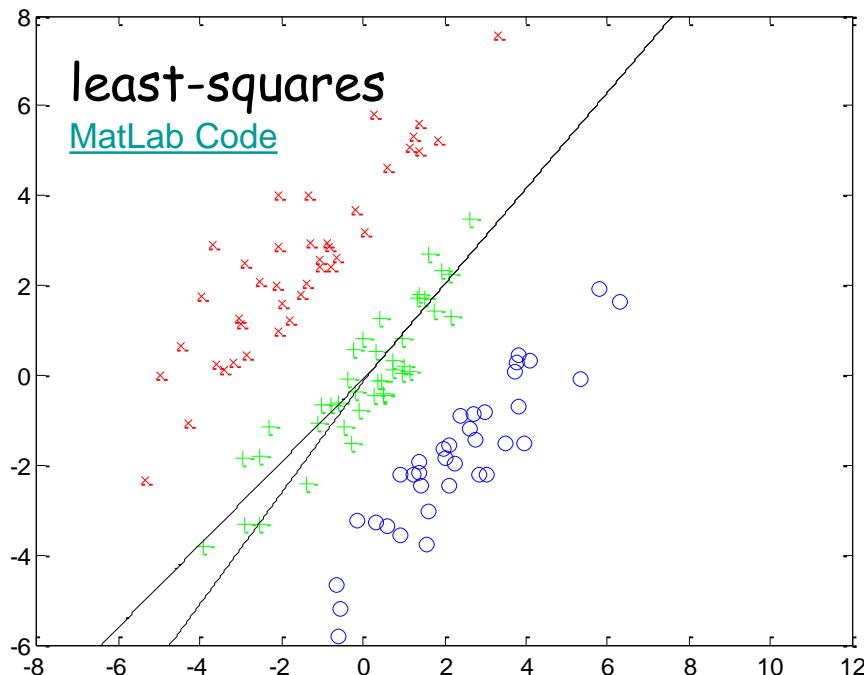
# Decision Boundaries for the K Class Case

- For the LS case, we show using the 3 activation functions the potential decision boundaries. A region is assigned to class  $j$  when  $y_j(x) > y_k(x)$  for all  $k \neq j$ . Verify the class assignments.



# Least Squares for Multiclass Classification

- In this example of 3 classes, logistic regression gives a good solution.
- The least-squares solution gives poor results. Least squares correspond to MLE under the assumption of a Gaussian conditional distribution, whereas binary target vectors have a distribution that is far from Gaussian.

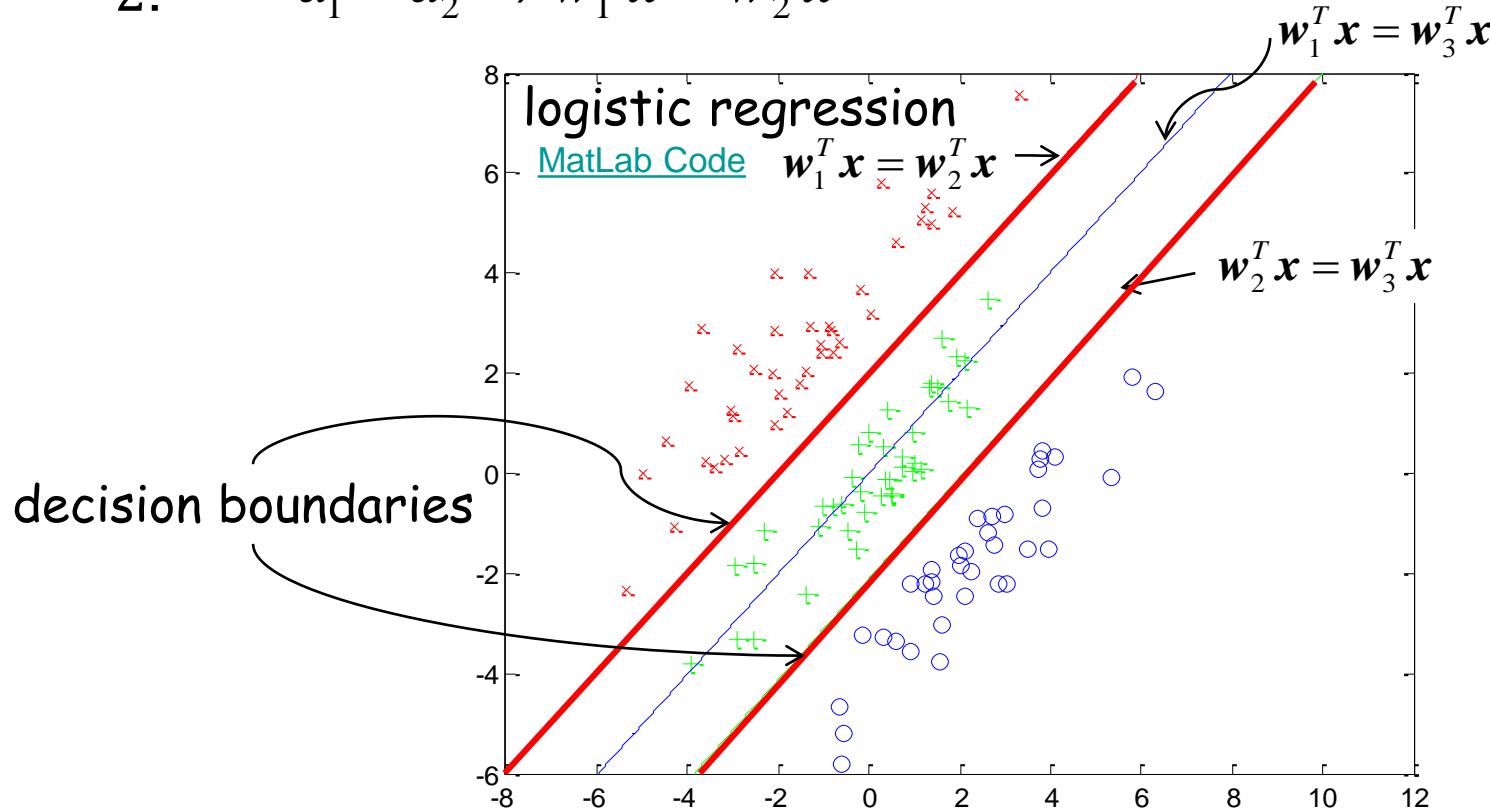


# Logistic Regression for Multiclass Classification

- For logistic regression we work with posterior class probabilities

$$p(C_k | \mathbf{x}) = \exp(a_k) / \sum_j \exp(a_j)$$

- At the decision boundaries (separating 2 classes), the corresponding posteriors are equal, e.g between classes 1 and 2:  $a_1 = a_2 \Rightarrow \mathbf{w}_1^T \mathbf{x} = \mathbf{w}_2^T \mathbf{x}$



---

# *Generative Vs. Discriminative Classifiers*

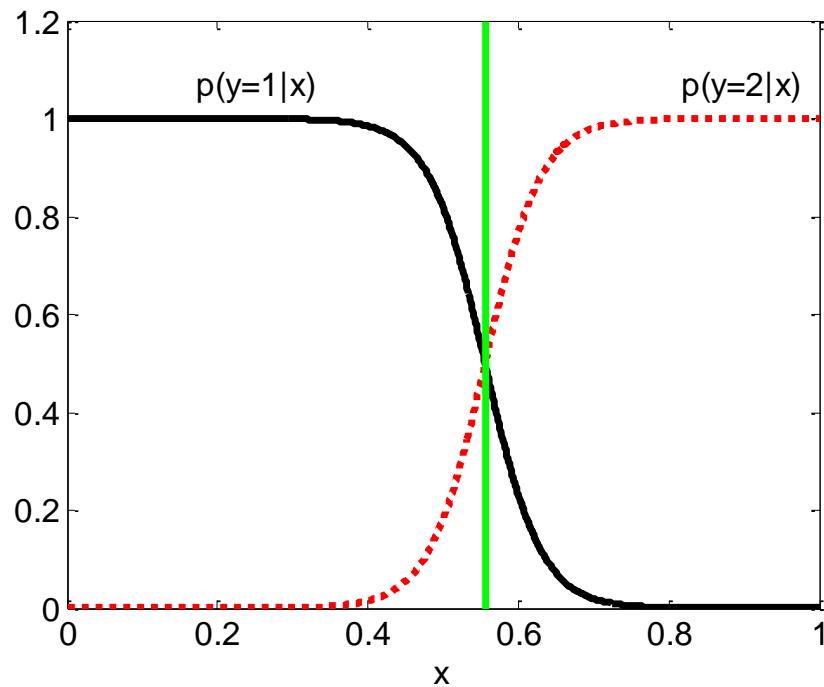
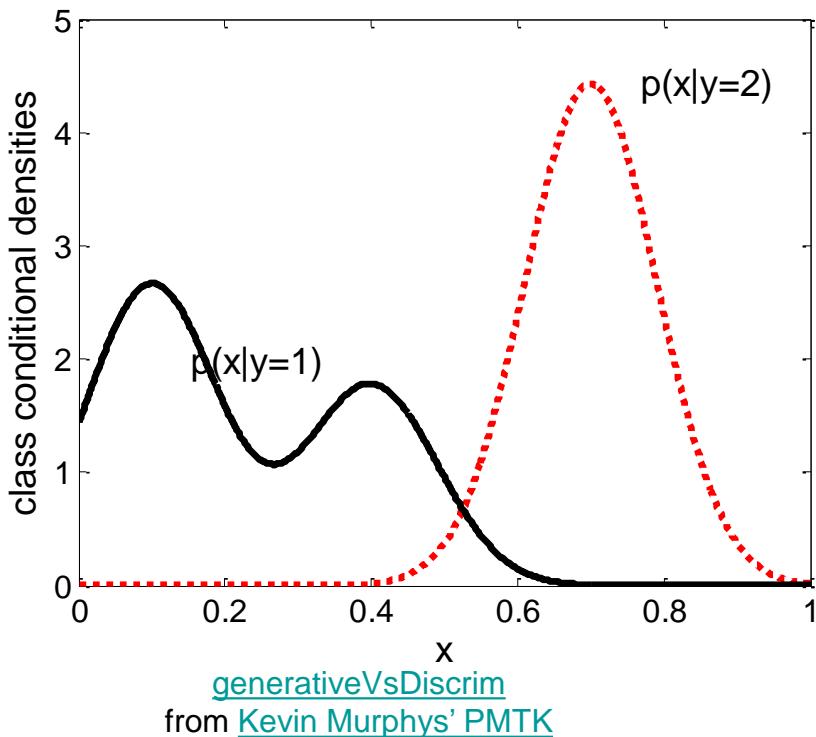
# Generative Versus Discriminative Classifiers

---

- *The assumptions made by GDA are much stronger than the assumptions made by logistic regression.*
  - When fitting a discriminative model, we usually maximize the conditional log likelihood  $\sum_{i=1}^N \log p(y_i | x_i, \theta)$ , whereas
  - when fitting a generative model, we usually maximize the joint log likelihood,  $\sum_{i=1}^N \log p(y_i, x_i | \theta)$ .
  - In general they give different results.
- When the Gaussian assumptions made by GDA are correct, the model will need less training data than logistic regression, but if the Gaussian assumptions are incorrect, logistic regression will do better.\*
- This is because discriminative models do not need to model the distribution of the features. This is illustrated next.
  - ❖ [Ng, A. Y. and M. I. Jordan \(2002\). On discriminative vs. generative classifiers: A comparison of logistic regression and naive bayes. In NIPS- 14.](#)

# Generative Versus Discriminative Classifiers

- The class conditional densities are complex (here  $p(x|y = 1)$  is a multimodal distribution), but the class posterior,  $p(y = c|x)$ , is a simple sigmoidal function, centered on the threshold value of 0.55.
- In general, *discriminative methods will be more accurate.*



# Generative Versus Discriminative Classifiers

---

- Easy to fit generative classifiers: e.g. in naive Bayes model and an LDA model by simple counting/averaging vs complex optimization in logistic regression.
- Adding new classes in a generative classifier is not a problem: we estimate the parameters of each class conditional density independently (not the case in discriminative models).
- In generative classifiers, we can deal easily with missing features: However, in a discriminative classifier, there is no principled solution to this problem, since the model is conditioned on  $x$ .

$$p(y=c | \mathbf{x}_{2:D}, \boldsymbol{\theta}) \propto p(y=c | \boldsymbol{\theta}) p(\mathbf{x}_{2:D} | y=c, \boldsymbol{\theta}) = p(y=c | \boldsymbol{\theta}) \sum_{x_1} p(x_1, \mathbf{x}_{2:D} | y=c, \boldsymbol{\theta})$$

- Salojarvi, J., K. Puolamaki, and S. Klaski (2005). [On discriminative joint density modeling](#). In *Proc. European Conf. on Machine Learning*
- Marlin, B. (2008). [Missing Data Problems in Machine Learning](#). Ph.D. thesis, U. Toronto
- Lasserre, J., C. Bishop, and T. Minka (2006). [Principled hybrids of generative and discriminative models](#). In *CVPR*.
- Liang, F., S. Mukherjee, and M. West (2007). [Understanding the use of unlabelled data in predictive modelling](#). *Statistical Science* 22, 189–205.
- Little., R. J. and D. B. Rubin (1987). [Statistical Analysis with Missing Data](#). New York: Wiley and Son.

# Generative Versus Discriminative Classifiers

---

$$p(y=c | \mathbf{x}_{2:D}, \boldsymbol{\theta}) \propto p(y=c | \boldsymbol{\theta}) p(\mathbf{x}_{2:D} | y=c, \boldsymbol{\theta}) = p(y=c | \boldsymbol{\theta}) \sum_{x_1} p(x_1, \mathbf{x}_{2:D} | y=c, \boldsymbol{\theta})$$

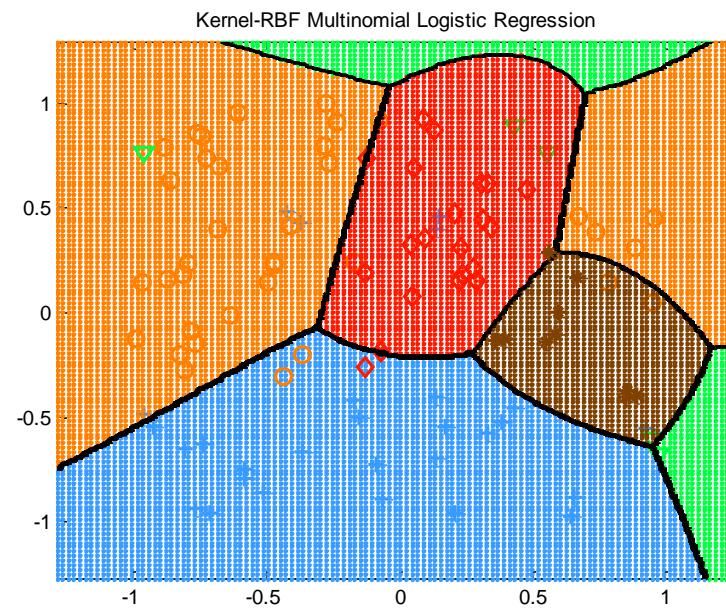
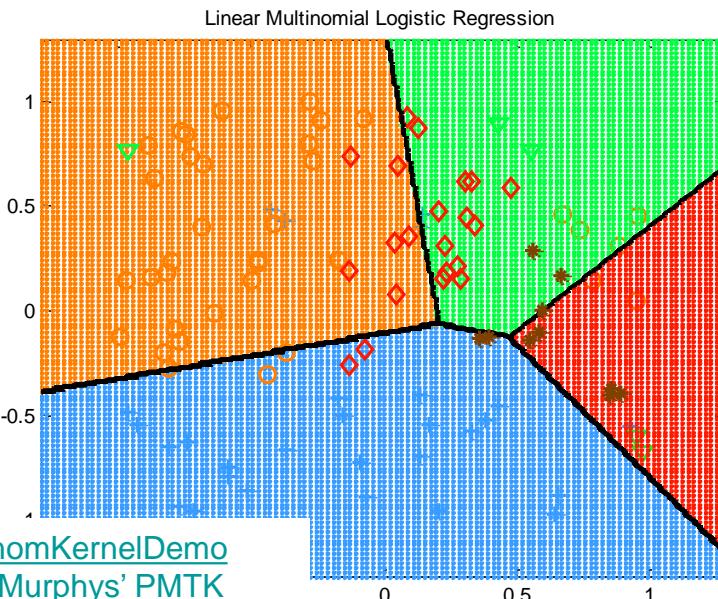
- This is further simplified if we employ the naïve Bayes assumption

$$\sum_{x_1} p(x_1, \mathbf{x}_{2:D} | y=c, \boldsymbol{\theta}) = \sum_{x_1} p(x_1 | y=c, \boldsymbol{\theta}_{1c}) \prod_{j=2}^D p(x_j | y=c, \boldsymbol{\theta}_{jc}) = \prod_{j=2}^D p(x_j | y=c, \boldsymbol{\theta}_{jc})$$

- In discriminative analysis, we can analytically marginalize missing values, e.g.  $p(\mathbf{x}_{2:D} | y=c, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}_{2:D} | \boldsymbol{\mu}_{c,2:D}, \boldsymbol{\Sigma}_{c,2:D,2:D})$
- Using unlabeled data in semi-supervised learning is easy with generative models: A much harder task in discriminative models.

# Generative Versus Discriminative Classifiers

- Generative models can generate inputs given the outputs: we can *infer probable inputs given the output*  $p(x|y)$ . This is not possible with a discriminative model. A generative model defines  $p(x, y)$ .
- *Discriminative methods allow preprocessing of the input:* replace  $x$  with  $\phi(x)$ . This cannot be done in generative models.
- *Discriminative Models – Better calibrated in terms of probability:* Independence assumptions in generative models (e.g. naive Bayes) lead to extreme posterior class probabilities (near 0 or 1).



# *Fisher's Linear Discriminant*

# Fisher's Linear Discriminant

---

- We now view a linear classification model in terms of **dimensionality reduction**. Consider two classes, and suppose we take the  $D$  –dimensional input  $x$  and project it to 1D:

$$y = \mathbf{w}^T \mathbf{x}$$

- If we place a threshold on  $y$  and classify  $y \geq -w_0$  as class  $C_1$ , and otherwise class  $C_2$ , then we obtain our standard linear classifier.
- The projection onto 1D leads to loss of information:
  - classes that are well separated in the original  $D$ -dimensional input space may become strongly overlapping in 1D.
- By adjusting  $w$ , can we select a projection that maximizes the class separation?

# Fisher's Linear Discriminant

---

- Consider a two-class problem in which there are  $N_1$  points of class  $C_1$  and  $N_2$  points of class  $C_2$ , so that the mean vectors of the two classes are given by

Class means:  $\mathbf{m}_1 = \frac{1}{N_1} \sum_{n \in C_1} \mathbf{x}_n, \quad \mathbf{m}_2 = \frac{1}{N_2} \sum_{n \in C_2} \mathbf{x}_n$

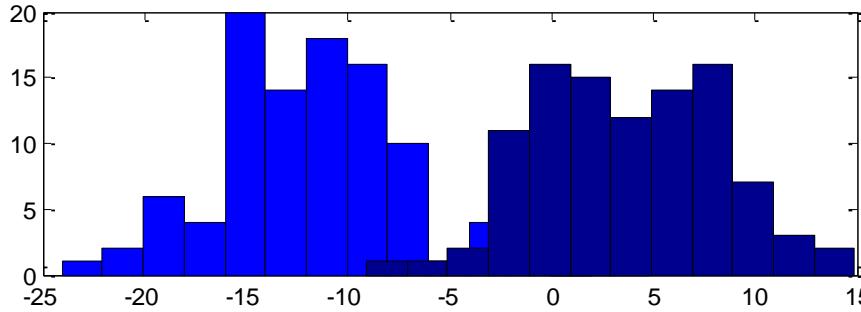
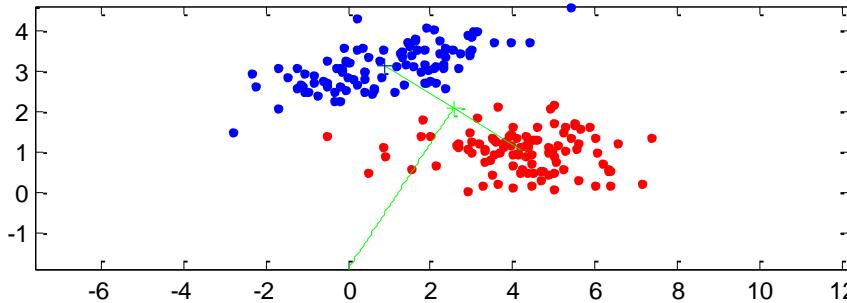
- A measure of the separation of the two classes when projected onto  $w$  is the separation of the projected means:

$$m_2 - m_1 = w^T (\mathbf{m}_2 - \mathbf{m}_1), \text{ where } m_i = w^T \mathbf{m}_i, i = 1, 2, \sum_i w_i^2 = 1$$

- The separation  $w^T (\mathbf{m}_2 - \mathbf{m}_1)$  can be made arbitrarily large by increasing  $|w|$ , so we constrain  $w$  as shown.
- Using a Lagrange multiplier to perform the constrained maximization we find that  $w \propto (\mathbf{m}_2 - \mathbf{m}_1)$ .

# Fisher's Linear Discriminant

- We show here 2 classes that are well separated in the  $(x_1, x_2)$  space but overlap when projected onto  $w \propto (m_2 - m_1)$ .

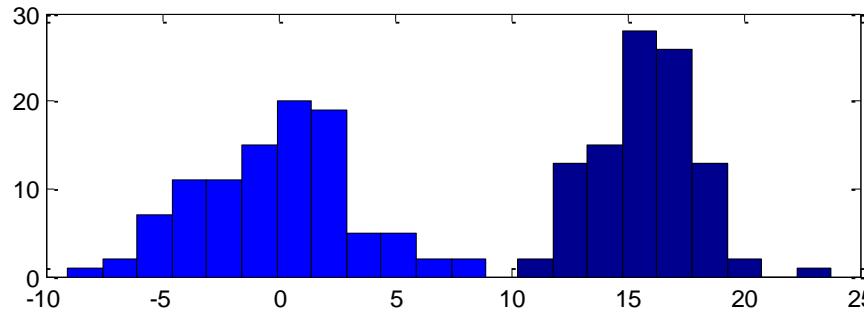
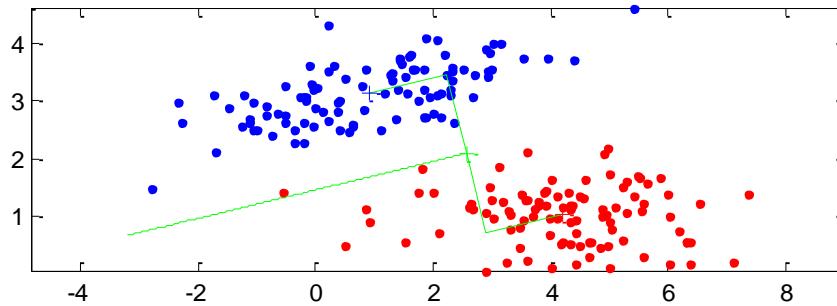


[MatLab Code](#)

- This is due to the *strongly non-diagonal covariances of the class distributions*.

# Fisher's Linear Discriminant

- We showed, in the last slide, two well separated classes in  $(x_1, x_2)$  space that overlap when projected onto  $w \propto (m_2 - m_1)$ .

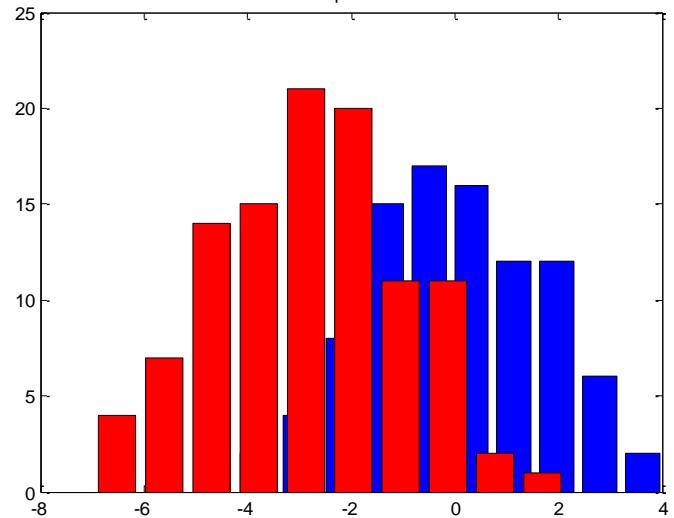
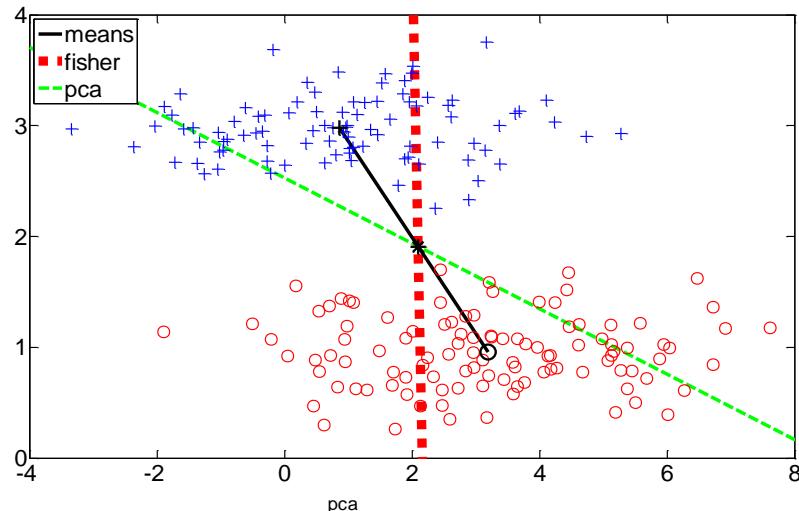


[MatLab Code](#)

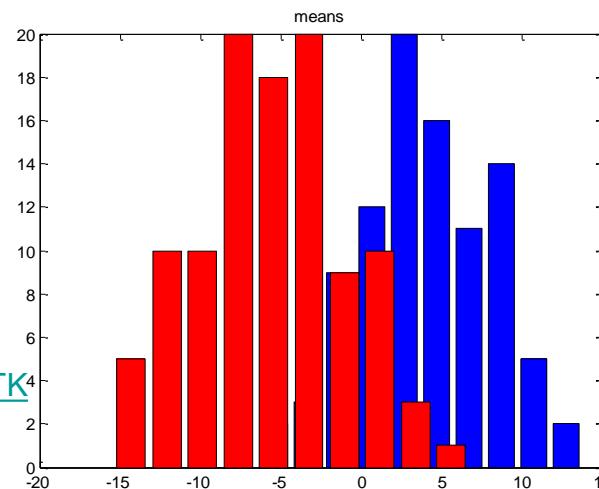
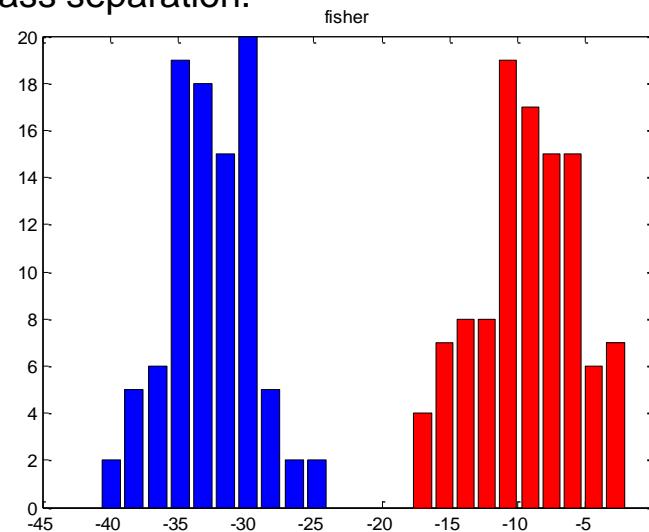
- Fisher proposed instead to maximize a function that gives a large separation between the projected class means while also giving a small variance within each class, thereby *minimizing the class overlap*.

# Fisher's Linear Discriminant

- (a) Comparison of class classification using 1<sup>st</sup> PCA vector, Fishers' linear discriminant vector and class-conditional means. (b) Projection of points onto Fisher's vector shows good class separation. (c) Projection of points onto PCA vector shows poor class separation.



[fisherLDAdemo](#)  
Kevin Murphys' PMTK



# Fisher's Linear Discriminant

---

- The within-class variance of the transformed data from class  $C_k$  is given by

$$s_k^2 = \sum_{n \in C_k} (y_n - m_k)^2, \text{ where } m_k = \mathbf{w}^T \mathbf{m}_k, k = 1, 2, \text{ and } y_n = \mathbf{w}^T \mathbf{x}_n$$

- We define the total within-class variance as  $s_1^2 + s_2^2$
- The Fisher criterion is then defined as the ratio of the between-class variance to the within-class variance for the transformed data as:

$$J(\mathbf{w}) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2} = \frac{\mathbf{w}^T S_B \mathbf{w}}{\mathbf{w}^T S_W \mathbf{w}}$$

where the between-class and within class covariances for the initial data are defined as:

$$S_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T, S_W = \sum_{n \in C_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T + \sum_{n \in C_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^T$$

# Fisher's Linear Discriminant

---

- Differentiating wrt  $\mathbf{w}$  gives:

$$\frac{d}{d\mathbf{w}} J(\mathbf{w}) = \frac{d}{d\mathbf{w}} \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} = 0 \Rightarrow (\mathbf{w}^T \mathbf{S}_B \mathbf{w}) \mathbf{S}_W \mathbf{w} = (\mathbf{w}^T \mathbf{S}_W \mathbf{w}) \mathbf{S}_B \mathbf{w}$$

- From  $\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T$ , we see that  $\mathbf{S}_B \mathbf{w}$  is in the direction of  $(\mathbf{m}_2 - \mathbf{m}_1)$ . Since we only care about the direction of  $\mathbf{w}$ , we drop the scalars  $(\mathbf{w}^T \mathbf{S}_B \mathbf{w})$ ,  $(\mathbf{w}^T \mathbf{S}_W \mathbf{w})$  and multiplying by  $\mathbf{S}_W^{-1}$  gives:

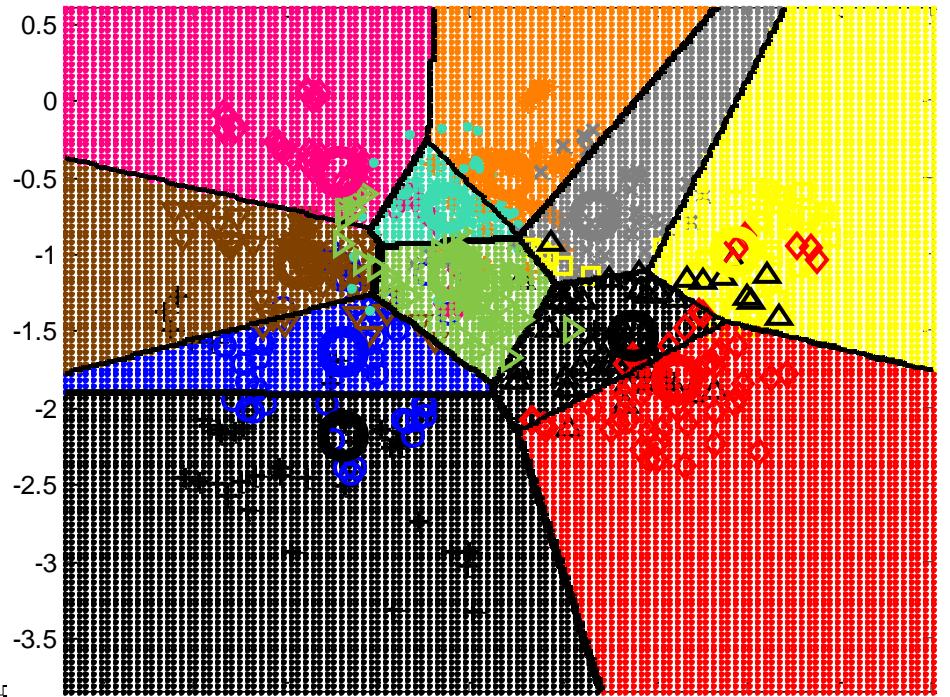
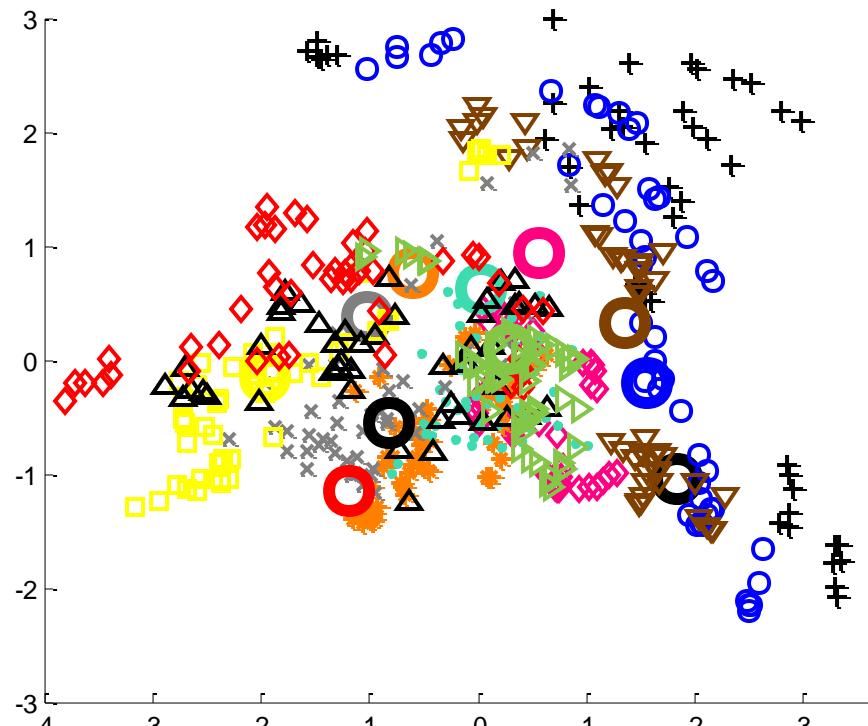
$$\mathbf{w} \propto \mathbf{S}_W^{-1}(\mathbf{m}_2 - \mathbf{m}_1)$$

$$\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T, \mathbf{S}_W = \sum_{n \in C_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T + \sum_{n \in C_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^T$$

- Note that if the within-class covariance  $\mathbf{S}_W$  is isotropic,  $\mathbf{S}_W \propto \mathbf{I}$ , then  $\mathbf{w} \propto (\mathbf{m}_2 - \mathbf{m}_1)$  as discussed earlier.

# Fisher's Linear Discriminant

- (a) PCA projection of vowel data (10-D) to 2d. (b) FLDA projection of vowel data to 2d. There is better class separation in the FLDA case.



[fisherDiscrimVowelDemo](#)  
from Kevin Murphy's PMTK

- Hastie, T., R. Tibshirani, and J. Friedman (2009). [The Elements of Statistical Learning](#). Springer. 2nd edition.

# Fisher's Linear Discriminant

---

$$\mathbf{w} \propto \mathbf{S}_W^{-1}(\mathbf{m}_2 - \mathbf{m}_1)$$

- The Fisher's linear discriminant is a choice of direction for projection of the data to 1D.
- The projected data can be used to construct a discriminant by choosing a threshold  $y_0$  so that we classify a new point belonging to  $C_1$  if  $y(\mathbf{x}) \geq y_0$  and to  $C_2$  otherwise.
- To choose the threshold, *we model the class-conditional densities  $p(y|C_k)$  using Gaussians and then use MLE to estimate their parameters.*
- Using the Gaussian approximations to the projected classes, *Decision theory (minimizing the misclassification rate) gives an expression for the optimal threshold.*
- Justification for the Gaussian assumption comes from the CLT noting that  $y = \mathbf{w}^T \mathbf{x}$  is sum of a set of random variables.

# Fisher's Discriminant vs. Least Squares

---

- The least-squares approach to determining a linear discriminant was based on making the model predictions as close as possible to the target values (using the 1-of-K coding).
- By contrast, the Fisher criterion was derived by requiring maximum class separation in the output space.
- For the two-class problem, the Fisher criterion can be obtained as a special case of least squares.
- Unfortunately we need to introduce a different target coding scheme (from the 1-of-K coding) for the least-squares solution for  $w$  to become equivalent to the Fisher solution.

# Fisher's Discriminant vs. Least Squares

- We take the targets for class  $C_1$  to be  $N/N_1$ , where  $N_1$  is the number of patterns in class  $C_1$ , and  $N$  is the total number of patterns. *This approximates the reciprocal of the prior probability for class  $C_1$ .*
- For class  $C_2$ , we take the targets to be  $-N/N_2$ , where  $N_2$  is the number of patterns in class  $C_2$ . With this choice:

$$\sum_{n=1}^N t_n = N_1 \frac{N}{N_1} + N_2 \left(-\frac{N}{N_2}\right) = 0$$

- The sum-of-squares error function is written as

$$E = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n)^2$$

- Taking the derivatives wrt  $w_0$  and  $\mathbf{w}$  gives:

$$\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n) = 0 \Rightarrow N\mathbf{w}^T \mathbf{m} + Nw_0 - \sum_{n=1}^N t_n = 0 \Rightarrow w_0 = -\mathbf{w}^T \mathbf{m}, \mathbf{m} = \frac{\sum_{n=1}^N \mathbf{x}_n}{N} = \frac{N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2}{N}$$
$$\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n) \mathbf{x}_n = 0 \Rightarrow \left( \mathbf{S}_w + \frac{N_1 N_2}{N} \mathbf{S}_B \right) \mathbf{w} = N(\mathbf{m}_1 - \mathbf{m}_2) \quad (\text{see proof on next slide})$$

# Fisher's Discriminant vs. Least Squares

$\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n) \mathbf{x}_n^T = 0$ , where  $w_0 = -\mathbf{w}^T \mathbf{m}$ ,  $\mathbf{m} = \frac{N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2}{N}$  and

$$\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T, \mathbf{S}_W = \sum_{n \in C_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T + \sum_{n \in C_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^T$$

□ The first Equ. can be simplified as:

$$\mathbf{w}^T \sum_{n \in C_1} \mathbf{x}_n \mathbf{x}_n^T + \mathbf{w}^T \sum_{n \in C_2} \mathbf{x}_n \mathbf{x}_n^T - \mathbf{w}^T \mathbf{m} (N \mathbf{m}^T) - \sum_{n=1}^N t_n \mathbf{x}_n^T = 0$$

□ Note from our choice of target values:

$$\sum_{n=1}^N t_n \mathbf{x}_n^T = \frac{N}{N_1} \sum_{n \in C_1} \mathbf{x}_n^T - \frac{N}{N_2} \sum_{n \in C_2} \mathbf{x}_n^T = N (\mathbf{m}_1^T - \mathbf{m}_2^T)$$

and also

$$\begin{aligned} \sum_{n \in C_1} \mathbf{x}_n \mathbf{x}_n^T &= \sum_{n \in C_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T + \mathbf{m}_1 \sum_{n \in C_1} \mathbf{x}_n^T + \sum_{n \in C_1} \mathbf{x}_n \mathbf{m}_1^T - N_1 \mathbf{m}_1 \mathbf{m}_1^T = \\ &= \sum_{n \in C_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T + N_1 \mathbf{m}_1 \mathbf{m}_1^T \end{aligned}$$

# Fisher's Discriminant vs. Least Squares

$$\mathbf{w}^T \sum_{n \in C_1} \mathbf{x}_n \mathbf{x}_n^T + \mathbf{w}^T \sum_{n \in C_2} \mathbf{x}_n \mathbf{x}_n^T - \mathbf{w}^T \mathbf{m} (N \mathbf{m}^T) - \sum_{n=1}^N t_n \mathbf{x}_n^T = 0$$

□ We now simplify the above equation as:

$$\begin{aligned} & \mathbf{w}^T \left\{ \sum_{n \in C_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T + N_1 \mathbf{m}_1 \mathbf{m}_1^T + \sum_{n \in C_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^T + N_2 \mathbf{m}_2 \mathbf{m}_2^T - \mathbf{m} (N \mathbf{m}^T) \right\} \\ & - N (\mathbf{m}_1^T - \mathbf{m}_2^T) = 0 \end{aligned}$$

or using the definition of  $S_W$

$$\begin{aligned} & \mathbf{w}^T \left\{ S_W + N_1 \mathbf{m}_1 \mathbf{m}_1^T + N_2 \mathbf{m}_2 \mathbf{m}_2^T - N \mathbf{m} \mathbf{m}^T \right\} = N (\mathbf{m}_1^T - \mathbf{m}_2^T) \Rightarrow \\ & \mathbf{w}^T \left\{ S_W + N_1 \mathbf{m}_1 \mathbf{m}_1^T + N_2 \mathbf{m}_2 \mathbf{m}_2^T - N \left( \frac{N_1}{N} \mathbf{m}_1 + \frac{N_2}{N} \mathbf{m}_2 \right) \left( \frac{N_1}{N} \mathbf{m}_1^T + \frac{N_2}{N} \mathbf{m}_2^T \right) \right\} = N (\mathbf{m}_1^T - \mathbf{m}_2^T) \Rightarrow \\ & \underbrace{\frac{N_1 N_2}{N} (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T}_{= \frac{N_1 N_2}{N} S_B} = N (\mathbf{m}_1^T - \mathbf{m}_2^T) \Rightarrow \\ & \left( S_w + \frac{N_1 N_2}{N} S_B \right) \mathbf{w} = N (\mathbf{m}_1 - \mathbf{m}_2) \end{aligned}$$

# Fisher's discriminant vs. Least Squares

---

$$\left( \mathbf{S}_w + \frac{N_1 N_2}{N} \mathbf{S}_B \right) \mathbf{w} = N(\mathbf{m}_1 - \mathbf{m}_2)$$

- Using  $\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T$ , we note that  $\mathbf{S}_B \mathbf{w}$  is in the direction of  $(\mathbf{m}_2 - \mathbf{m}_1)$ . Thus we can write

$$\mathbf{w} \propto \mathbf{S}_w^{-1} (\mathbf{m}_2 - \mathbf{m}_1)$$

- This is the same result obtained with the Fisher criterion!
- We also found the bias  $w_0 = -\mathbf{w}^T \mathbf{m}$ . This tells us that a new vector  $x$  should be classified as belonging to  $C_1$  if

$$y(x) = w_0 + \mathbf{w}^T x = -\mathbf{w}^T \mathbf{m} + \mathbf{w}^T x = \mathbf{w}^T (x - \mathbf{m}) > 0$$

and class  $C_2$  otherwise.

# Fisher's Discriminant for Multiple Classes

---

- We now consider the generalization of the Fisher discriminant to  $K > 2$  classes.
- Assume the dimensionality of the input space is  $D > K$ .
- We introduce  $D' > 1$  linear ‘features’  $y_k = \mathbf{w}_k^T \mathbf{x}, k = 1, \dots, D'$ . We group these feature values in a vector  $\mathbf{y}$ .
- Similarly, the weight vectors  $\{\mathbf{w}_k\}$  are taken as the columns of a matrix  $\mathbf{W}$ , so that

$$\mathbf{y} = \mathbf{W}^T \mathbf{x}$$

- Note that here we are not including any bias parameters in the definition of  $\mathbf{y}$ .
- The within-class covariance matrix for  $K$  classes is:

$$\mathbf{S}_W = \sum_{k=1}^K \mathbf{S}_k, \mathbf{S}_k = \sum_{n \in C_k} (\mathbf{x}_n - \mathbf{m}_k)(\mathbf{x}_n - \mathbf{m}_k)^T, \mathbf{m}_k = \frac{1}{N_k} \sum_{n \in C_k} \mathbf{x}_n$$

# Fisher's Discriminant for Multiple Classes

- We need to also find the between-class covariance matrix. We work here in the  $x$  –space. Start with the total covariance:

$$S_T = \sum_{n=1}^N (\mathbf{x}_n - \mathbf{m})(\mathbf{x}_n - \mathbf{m})^T, \quad \mathbf{m} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \frac{1}{N} \sum_{k=1}^K N_k \mathbf{m}_k, \quad \sum_{k=1}^K N_k = N$$

- The total covariance matrix can be decomposed into the sum of the **within-class covariance matrix** plus an additional matrix  $S_B$ , which we identify as a measure of **the between-class covariance**:

$$S_T = S_W + S_B, \quad S_W = \sum_{k=1}^K S_k, \quad S_k = \sum_{n \in C_k} (\mathbf{x}_n - \mathbf{m}_k)(\mathbf{x}_n - \mathbf{m}_k)^T, \quad \mathbf{m}_k = \frac{1}{N_k} \sum_{n \in C_k} \mathbf{x}_n$$

- From the above two Eqs we can show that:

$$S_B = \sum_k \sum_{n \in C_k} (\mathbf{x}_n - \mathbf{m})(\mathbf{x}_n - \mathbf{m})^T - \sum_k \sum_{n \in C_k} (\mathbf{x}_n - \mathbf{m}_k)(\mathbf{x}_n - \mathbf{m}_k)^T = \sum_{k=1}^K N_k (\mathbf{m}_k - \mathbf{m})(\mathbf{m}_k - \mathbf{m})^T$$

# Fisher's Discriminant for Multiple Classes

- We can define similar covariance matrices in the  $D'$  dimensional  $\mathbf{y}$ -space:

$$\mathbf{S}_W = \sum_{k=1}^K \sum_{n \in C_k} (\mathbf{y}_n - \boldsymbol{\mu}_k)(\mathbf{y}_n - \boldsymbol{\mu}_k)^T, \quad \boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n \in C_k} \mathbf{y}_n$$

$$\mathbf{S}_B = \sum_{k=1}^K N_k (\boldsymbol{\mu}_k - \boldsymbol{\mu})(\boldsymbol{\mu}_k - \boldsymbol{\mu})^T, \quad \boldsymbol{\mu} = \frac{1}{N} \sum_{k=1}^K N_k \boldsymbol{\mu}_k$$

- We wish to maximize the between-class covariance while minimizing the within-class covariance. There are many choices of (scalar) criteria:

$$J(\mathbf{W}) = \text{Tr} \left\{ \mathbf{S}_W^{-1} \mathbf{S}_B \right\} = \text{Tr} \left\{ (\mathbf{W}^T \mathbf{S}_W \mathbf{W})^{-1} (\mathbf{W}^T \mathbf{S}_B \mathbf{W}) \right\}$$

- Maximization is straightforward. The weight values are determined as  $\mathbf{W} = \mathbf{S}_W^{-1/2} \mathbf{U}$  where  $\mathbf{U}$  are the  $D'$  leading eigenvectors of  $\mathbf{S}_W^{-1/2} \mathbf{S}_B \mathbf{S}_W^{-1/2}$  ( $D'$  = number of linear features).

# Fisher's Discriminant for Multiple Classes

---

- Note  $S_B = \sum_{k=1}^K N_k (\mathbf{m}_k - \mathbf{m})(\mathbf{m}_k - \mathbf{m})^T$  is composed of the sum of  $K$  matrices, each of which is an outer product of two vectors and therefore of rank 1.
- However, only  $(K - 1)$  of these matrices are independent since  $\mathbf{m} = \sum_{k=1}^K N_k \mathbf{m}_k / N$ . Thus,  $\mathbf{S}_B$  has rank at most  $(K - 1)$  i.e. has at most  $(K - 1)$  nonzero eigenvalues.
- Projecting onto the  $(K - 1)$  –dim subspace spanned by the eigenvectors of  $\mathbf{S}_B$  does not change the value of  $J(\mathbf{W})$ .
  - *Severe Limitation of FLDA: We can find at most  $L = K - 1$  linear features regardless of the dimensionality of  $x$ .*
  - Fukunaga, K. (1990). *Introduction to Statistical Pattern Recognition* (Second ed.). Academic Press.

# Probabilistic Interpretation of FLDA

---

- A probabilistic interpretation of FLDA was proposed.<sup>1,2</sup>
- In the **heteroscedastic LDA (HLDA)**, let  $\mathbf{W}$  be a  $D \times D$  invertible matrix and let  $\mathbf{z}_i = \mathbf{W}\mathbf{x}_i$  a transformed version of the data.
- We *fit full covariance Gaussians to the transformed data, one per class*, but with the constraint that *only the first  $L$  components will be class-specific*.
- The remaining  $H = D - L$  components are shared across classes, and will thus not be discriminative.

$$p(\mathbf{z}_i | \boldsymbol{\theta}, y_i = c) = \mathcal{N}(\mathbf{z}_i | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c), \boldsymbol{\mu}_c = \begin{pmatrix} \mathbf{m}_c \\ \mathbf{m}_0 \end{pmatrix}, \boldsymbol{\Sigma}_c = \begin{pmatrix} \mathbf{S}_c & 0 \\ 0 & \mathbf{S}_0 \end{pmatrix}$$

1. Kumar, N. and A. Andreo (1998). [Heteroscedastic discriminant analysis and reduced rank HMMs for improved speech recognition](#). *Speech Communication* 26, 283–297.
2. Zhou, H., D. Karakos, S. Khudanpur, A. Andreou, and C. Priebe (2009). [On Projections of Gaussian Distributions using Maximum Likelihood Criteria](#). In *Proc. of the Workshop on Information Theory and its Applications*.
3. Gales, M. (2002). [Maximum likelihood multiple subspace projections for hidden Markov models](#). *IEEE Trans. on Speech and Audio Processing* 10(2), 37–47.
4. Gales, M. J. F. (1999). [Semi-tied covariance matrices for hidden Markov models](#). *IEEE Trans. on Speech and Audio Processing* 7(3), 272–281.

# Probabilistic Interpretation of FLDA

---

$$p(z_i | \theta, y_i = c) = \mathcal{N}(z_i | \mu_c, \Sigma_c), \quad \mu_c = \begin{pmatrix} \mathbf{m}_c \\ \mathbf{m}_0 \end{pmatrix}, \quad \Sigma_c = \begin{pmatrix} \mathbf{S}_c & 0 \\ 0 & \mathbf{S}_0 \end{pmatrix}$$

- Here  $\mathbf{m}_0$  is the shared  $H$  dimensional mean and  $\mathbf{S}_0$  is the shared  $H \times H$  covariace.
- The pdf of the original (untransformed) data is then:

$$p(\mathbf{x}_i | y_i = c, \mathbf{W}, \boldsymbol{\theta}) = |\mathbf{W}| \mathcal{N}(\mathbf{W}\mathbf{x}_i | \mu_c, \Sigma_c) = |\mathbf{W}| \mathcal{N}(\mathbf{W}_L \mathbf{x}_i | \mathbf{m}_c, \mathbf{S}_c) \mathcal{N}(\mathbf{W}_H \mathbf{x}_i | \mathbf{m}_0, \mathbf{S}_0), \quad \mathbf{W} = \begin{pmatrix} \mathbf{W}_L \\ \mathbf{W}_H \end{pmatrix}$$

- For fixed  $\mathbf{W}$ , it is easy to derive the MLE for  $\boldsymbol{\theta}$ . One can then optimize  $\mathbf{W}$  using gradient methods.

# Probabilistic Interpretation of FLDA

---

$$p(z_i | \theta, y_i = c) = \mathcal{N}(z_i | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c), \boldsymbol{\mu}_c = \begin{pmatrix} \mathbf{m}_c \\ \mathbf{m}_0 \end{pmatrix}, \boldsymbol{\Sigma}_c = \begin{pmatrix} S_c & 0 \\ 0 & S_0 \end{pmatrix}$$

- When  $\boldsymbol{\Sigma}_c$  are diagonal, there is a closed-form for  $\mathbf{W}$ .<sup>2</sup>
- In the special case the  $\boldsymbol{\Sigma}_c$  are all equal, we recover classical LDA.<sup>4</sup>
- HLDA outperforms LDA if the class covariances are not equal within the discriminative subspace (i.e., **if the assumption that  $\boldsymbol{\Sigma}_c$  is independent of  $c$  is a poor assumption**). This was demonstrated on speech recognition.<sup>3</sup>
- One can allow each class to use its own projection matrix (**multiple LDA**).<sup>1</sup>

1. Gales, M. (2002). [Maximum likelihood multiple subspace projections for hidden Markov models](#). *IEEE Trans. on Speech and Audio Processing* 10(2), 37–47
2. Gales, M. J. F. (1999). [Semi-tied covariance matrices for hidden Markov models](#). *IEEE Trans. on Speech and Audio Processing* 7 (3), 272–281.
3. Kumar, N. and A. Andreou (1998). [Heteroscedastic discriminant analysis and reduced rank HMMs for improved speech recognition](#). *Speech Communication* 26, 283–297.
4. Zhou, H., D. Karakos, S. Khudanpur, A. Andreou, and C. Priebe (2009). [On Projections of Gaussian Distributions using Maximum Likelihood Criteria](#). In *Proc. of the Workshop on Information Theory and its Applications*.

---

# *Probabilistic Generative Models*

# Probabilistic Generative Models

---

- Models with linear decision boundaries arise from simple assumptions about the probabilistic distribution of the data.
- We here adopt **a generative approach**: we model the class-conditional densities  $p(\mathbf{x}|C_k)$ , as well as the class priors  $p(C_k)$ , and then use these to compute posterior probabilities  $p(C_k|\mathbf{x})$  through Bayes' theorem.
- Let us consider two classes:

$$p(C_1 | \mathbf{x}) = \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_1)p(C_1) + p(\mathbf{x} | C_2)p(C_2)} = \frac{1}{1 + e^{-a}} = \sigma(a)$$

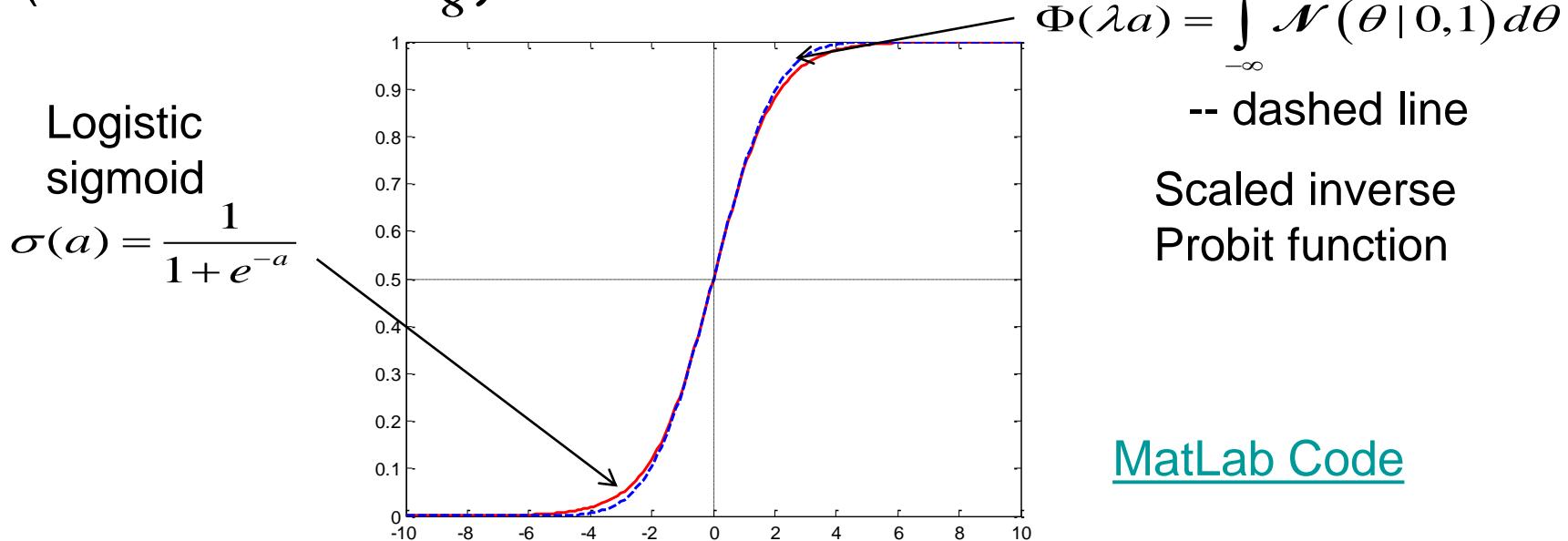
where:

$$a = \ln \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_2)p(C_2)} \quad \text{and} \quad \sigma(a) = \frac{1}{1 + e^{-a}}$$

Logistic sigmoid

# Logistic Sigmoid (S-shaped)

- We can approximate the logistic sigmoid function with the scaled inverse probit function.
- The derivatives of the two curves are the same for  $a = 0$  (selection of  $\lambda^2 = \frac{\pi^2}{8}$ ).



[MatLab Code](#)

- The logistic sigmoid function satisfies:  
$$\sigma(-a) = 1 - \sigma(a)$$
- The *inverse of  $\sigma(a)$  is the logit function:*  
$$a = \ln \frac{\sigma}{1 - \sigma} = \ln \frac{p(C_1 | \mathbf{x})}{p(C_2 | \mathbf{x})}$$

# Probabilistic Generative Models

---

- The appearance of the logistic sigmoid is rather arbitrary up to this point.

$$p(C_1 | \mathbf{x}) = \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_1)p(C_1) + p(\mathbf{x} | C_2)p(C_2)} = \frac{1}{1+e^{-a}} = \sigma(a)$$

- This will be useful provided  $a(\mathbf{x}) = \ln \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_2)p(C_2)}$  takes a simple functional form.
- We will consider problems with  $a(\mathbf{x})$  being a linear function of  $\mathbf{x}$ .
- In this case, the posterior probability is governed by a generalized linear model.

# Probabilistic Generative Models $K > 2$

---

- For the case of  $K > 2$  classes, we have

$$p(C_k | \mathbf{x}) = \frac{p(\mathbf{x} | C_k)p(C_k)}{\sum_j p(\mathbf{x} | C_j)p(C_j)} = \frac{e^{a_k}}{\sum_j e^{a_j}}, \quad a_k = \ln(p(\mathbf{x} | C_k)p(C_k))$$

- This is known as the normalized exponential and can be regarded as a multiclass generalization of the logistic sigmoid.
- Here we define:

$$a_k = \ln(p(\mathbf{x} | C_k)p(C_k))$$

- The normalized exponential is known as the softmax function.
- It can be seen as a smoothed version of the max function because, if  $a_k \geq a_j$  for all  $j \neq k$ , then  $p(C_k | \mathbf{x}) \cong 1$ , and  $p(C_j | \mathbf{x}) \cong 0$ .

# Gaussian Class Conditionals

- Assume that the class-conditional densities are Gaussians and all classes share the same covariance matrix:

$$p(\mathbf{x} | C_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\}$$

- We consider the case of two classes:

$$p(C_1 | \mathbf{x}) = \sigma(a), a(\mathbf{x}) = \ln \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_2)p(C_2)}$$

where:

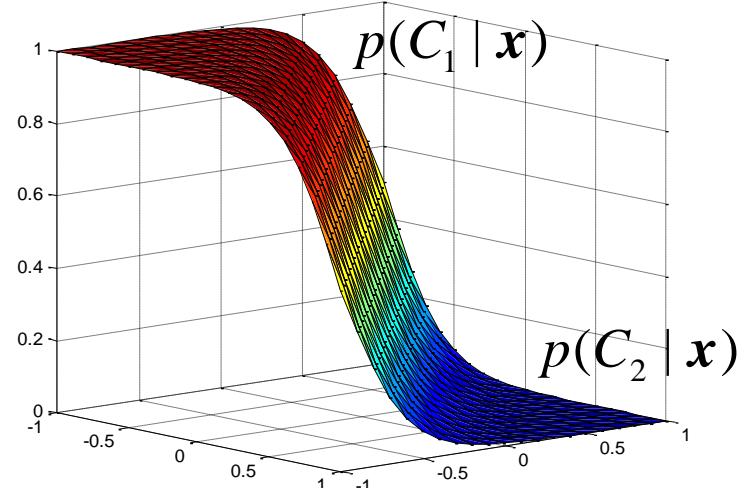
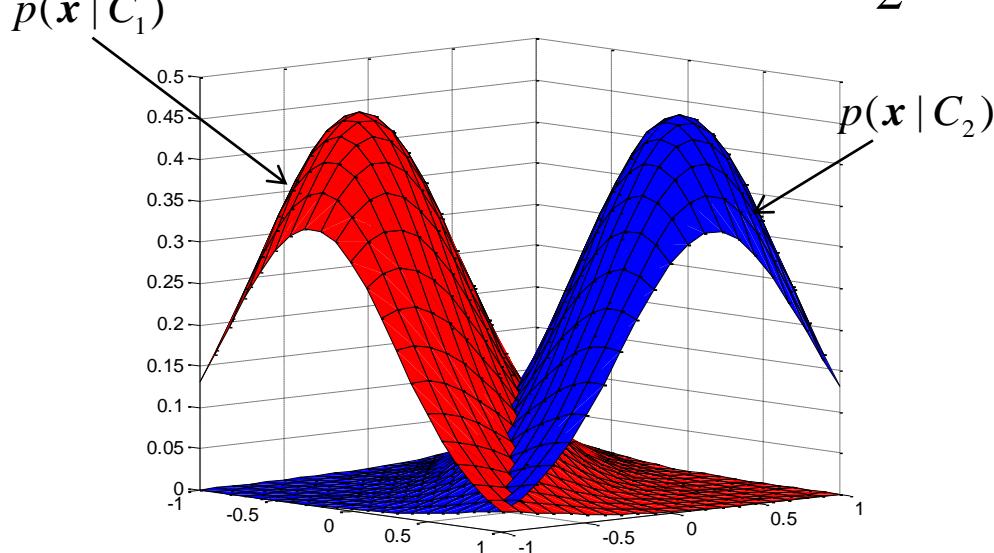
$$\begin{aligned} a(\mathbf{x}) &= -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_2)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) + \ln \frac{p(C_1)}{p(C_2)} = \\ &= \underbrace{\left( \boldsymbol{\mu}_1^T \Sigma^{-1} - \boldsymbol{\mu}_2^T \Sigma^{-1} \right) \mathbf{x}}_{w^T} - \underbrace{\frac{1}{2} \boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_2^T \Sigma^{-1} \boldsymbol{\mu}_2 + \ln \frac{p(C_1)}{p(C_2)}}_{w_0} \end{aligned}$$

# Gaussian Class Conditionals and Posteriors

- We thus obtain a linear function of  $x$ :

$$p(C_1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

$$\mathbf{w} = \Sigma^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2), w_0 = -\frac{1}{2} \boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_2^T \Sigma^{-1} \boldsymbol{\mu}_2 + \ln \frac{p(C_1)}{p(C_2)}$$



[MatLab Code](#)

- Class conditional densities (left) and corresponding posteriors (right).
- On the right, the surface is colored using a proportion of red ink given by  $p(C_1 | \mathbf{x})$  and a proportion of blue ink given by  $p(C_2 | \mathbf{x}) = 1 - p(C_1 | \mathbf{x})$

# Gaussian Class Posteriors

---

$$p(C_1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

- *The decision boundaries correspond to surfaces along which the posterior probabilities  $p(C_k | \mathbf{x})$  are constant.*
- They are linear in input space.
- The prior probabilities  $p(C_k)$  enter only through the bias parameter  $w_0$

$$w_0 = -\frac{1}{2} \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2 + \ln \frac{p(C_1)}{p(C_2)}$$

so *changing  $p(C_k)$  only leads to parallel shifts of the decision boundaries.*

# Probabilistic Generative Models $K > 2$

---

- For the case of  $K > 2$  classes, we have

$$p(C_k | \mathbf{x}) = \frac{p(\mathbf{x} | C_k)p(C_k)}{\sum_j p(\mathbf{x} | C_j)p(C_j)} = \frac{e^{a_k}}{\sum_j e^{a_j}} \quad \text{where } a_k = \ln(p(\mathbf{x} | C_k)p(C_k))$$

- The resulting decision boundaries, corresponding to the minimum misclassification rate, occur when the two largest posterior probabilities are equal, and so are defined by linear functions of  $\mathbf{x}$ , leading to a generalized linear model.
- For this case, we can see immediately that:

$$\begin{aligned} a_k(\mathbf{x}) &= \ln(p(\mathbf{x} | C_k)p(C_k)) = -\frac{1}{2}\mathbf{x}^T \Sigma^{-1} \mathbf{x} + \mathbf{w}_k^T \mathbf{x} + w_{k0} \\ \mathbf{w}_k &= \Sigma^{-1} \boldsymbol{\mu}_k, \quad w_{k0} = -\frac{1}{2} \boldsymbol{\mu}_k^T \Sigma^{-1} \boldsymbol{\mu}_k + \ln p(C_k) + \ln \frac{1}{(2\pi)^{D/2}} - \frac{1}{2} \ln |\Sigma| \end{aligned}$$

# Probabilistic Generative Models K>2

---

- Note that we can now re-write  $p(C_k | \mathbf{x})$  by cancelling out terms:

$$p(C_k | \mathbf{x}) = \frac{p(\mathbf{x} | C_k) p(C_k)}{\sum_j p(\mathbf{x} | C_j) p(C_j)} = \frac{e^{a_k}}{\sum_j e^{a_j}} \quad \text{where} \quad a_k \equiv \mathbf{w}_k^T \mathbf{x} + w_{k0}$$

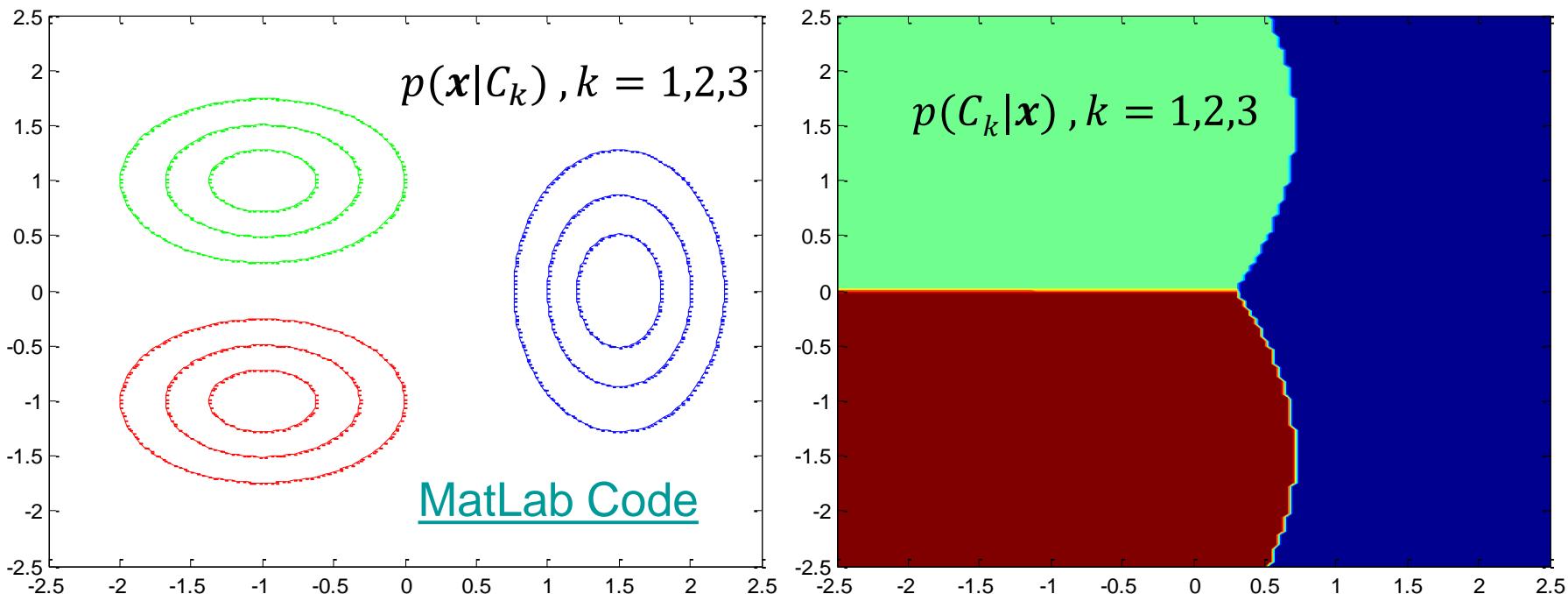
where:

$$\mathbf{w}_k = \Sigma^{-1} \boldsymbol{\mu}_k, w_{k0} = -\frac{1}{2} \boldsymbol{\mu}_k^T \Sigma^{-1} \boldsymbol{\mu}_k + \ln p(C_k)$$

- $a_k(\mathbf{x})$  are linear functions of  $\mathbf{x}$  due to the cancellation of the quadratic terms with the shared  $\Sigma$ .

# Probabilistic Generative Models $K > 2$

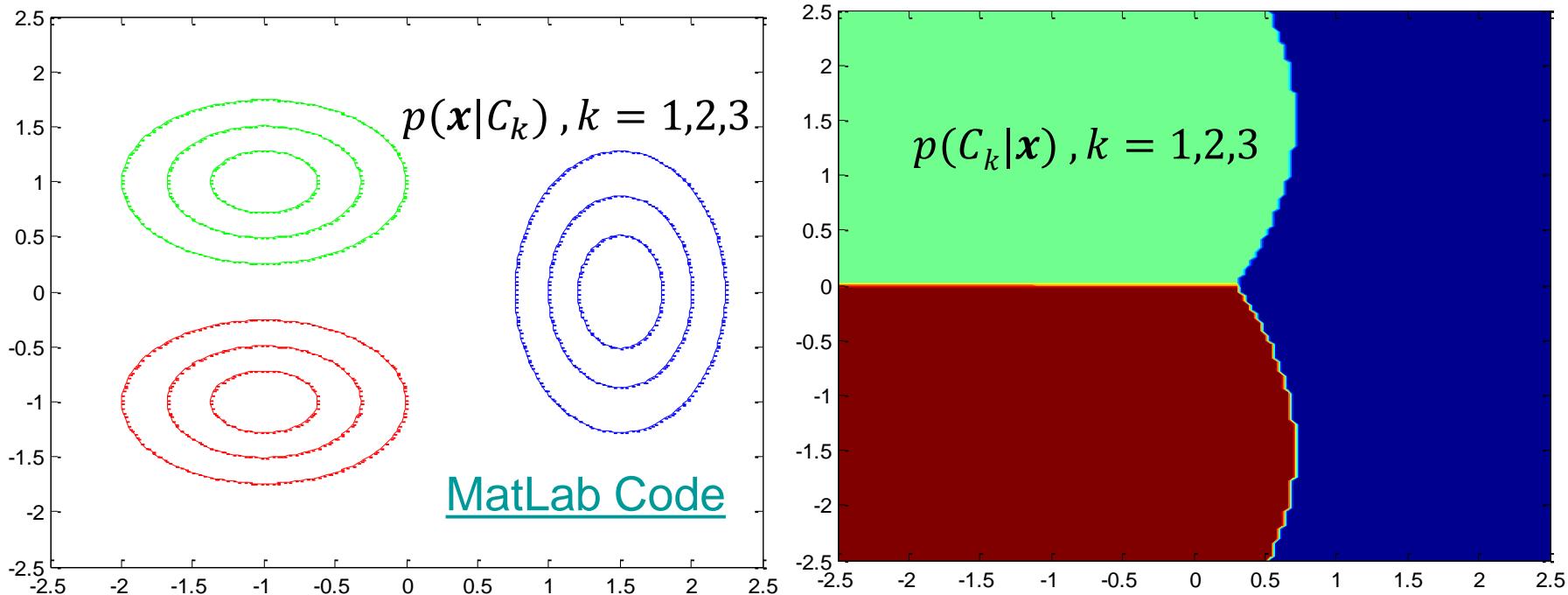
- If we relax the assumption of a shared covariance matrix and allow each  $p(x|C_k)$  to have its own  $\Sigma_k$ , then we obtain quadratic functions of  $x$  and a quadratic discriminant.



- The left-hand plot shows  $p(x|C_k), k = 1,2,3$  each a Gaussian distribution. The red & green classes have the same  $\Sigma$ .

# Probabilistic Generative Models $K > 2$

- The right-hand plot shows the posteriors  $p(C_k|x)$ ,  $k = 1,2,3$ , in which the three colors represent the posterior probabilities for the respective 3 classes (scaled).
- Note that for *the two classes with the same covariance, the decision boundary is linear while for the other 2 is quadratic.*



# Maximum Likelihood Solution

---

- Once we have a parametric form for  $p(\mathbf{x}|C_k)$ , we can compute the parameters, together with the prior class probabilities  $p(C_k)$ , using MLE.
- This requires a training data set comprising of  $\mathbf{x}$  and their corresponding class labels.
- Consider the case of two classes each having a Gaussian class-conditional density with a shared covariance matrix.
- Suppose we have a data set  $\{\mathbf{x}_n, t_n\}$  where  $n = 1, \dots, N$ . Here  $t_n = 1$  denotes class  $C_1$  and  $t_n = 0$  denotes  $C_2$ .
- We denote the prior class probabilities  $p(C_1) = \pi$ ,  $p(C_2) = 1 - \pi$ .

# Maximum Likelihood Solution

- For a data point  $x_n$  from class  $C_1$ , we have  $t_n = 1$  and hence

$$p(x_n, C_1) = p(C_1)p(x_n | C_1) = \pi \mathcal{N}(x_n | \mu_1, \Sigma)$$

- Similarly for class  $C_2$ , we have  $t_n = 0$

$$p(x_n, C_2) = p(C_2)p(x_n | C_2) = (1 - \pi) \mathcal{N}(x_n | \mu_2, \Sigma)$$

- Thus *the likelihood function is given by*

$$p(t, X / \pi, \mu_1, \mu_2, \Sigma) = \prod_{n=1}^N [\pi \mathcal{N}(x_n | \mu_1, \Sigma)]^{t_n} [(1 - \pi) \mathcal{N}(x_n | \mu_2, \Sigma)]^{1-t_n}$$

where  $t = (t_1, \dots, t_N)^T$ .

- We maximize the log likelihood function.

# Maximum Likelihood Solution

$$p(\mathbf{t}, \mathbf{X} / \pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = \prod_{n=1}^N [\pi \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma})]^{t_n} [(1-\pi) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma})]^{1-t_n}$$

- Consider first the maximization with respect to  $\pi$ . The terms of interest in the log-likelihood are:

$$\sum_{n=1}^N \{ t_n \ln \pi + (1-t_n) \ln(1-\pi) \}$$

- Taking derivative wrt  $\pi$  and setting it equal to zero:

$$\pi = \frac{1}{N} \sum_{n=1}^N t_n = \frac{N_1}{N_1 + N_2}$$

where  $N_i$  the total number of data points in class  $C_i$ .

- As expected, the maximum likelihood estimate for  $\pi$  is simply the fraction of points in class  $C_1$ .
- This result is easily generalized to the multiclass case (see next) where the MLE of the prior probability associated with class  $C_k$  is given by the fraction of the training set points  $N_k$ .

# Maximum Likelihood Solution: Multiclass

$$p(\{\phi_n, t_n\} / \{\pi_k\}) = \prod_{n=1}^N \prod_{k=1}^K [\pi_k p(\phi_n | C_k)]^{t_{nk}}$$

- For the multiclass case, the log likelihood takes the form:

$$\ln p(\{\phi_n, t_n\} / \{\pi_k\}) = \sum_{n=1}^N \sum_{k=1}^K t_{nk} (\ln \pi_k + \ln p(\phi_n | C_k)), \quad \sum_{k=1}^K \pi_k = 1$$

- To maximize the log likelihood, we consider:

$$\sum_{n=1}^N \sum_{k=1}^K t_{nk} (\ln \pi_k + \ln p(\phi_n | C_k)) + \lambda \left( \sum_{k=1}^K \pi_k - 1 \right)$$

- Taking the derivative wrt  $\pi_k$  we arrive at the same result as for two classes:

$$\sum_{n=1}^N \frac{t_{nk}}{\pi_k} + \lambda = 0 \Rightarrow \pi_k \lambda = -\sum_{n=1}^N t_{nk} = -N_k \Rightarrow \begin{cases} \sum_{k=1}^K \pi_k \lambda = -N \\ \pi_k \lambda = -N_k \end{cases} \Rightarrow \begin{cases} \lambda = -N \\ \pi_k \lambda = -N_k \end{cases} \Rightarrow \pi_k = \frac{N_k}{N}$$

# Maximum Likelihood Solution

$$p(\mathbf{t}, \mathbf{X} | \pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = \prod_{n=1}^N [\pi \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma})]^{t_n} [(1-\pi) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma})]^{1-t_n}$$

- Consider now the maximization with respect to  $\boldsymbol{\mu}_1$ . The terms of interest in the log-likelihood are:

$$\sum_{n=1}^N t_n \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}) = -\frac{1}{2} \sum_{n=1}^N t_n (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_1) + const.$$

- Taking derivative wrt  $\boldsymbol{\mu}_1$  and setting it equal to zero:

$$\boldsymbol{\mu}_1 = \frac{1}{N_1} \sum_{n=1}^N t_n \mathbf{x}_n, \text{ and similarly } \boldsymbol{\mu}_2 = \frac{1}{N_2} \sum_{n=1}^N (1-t_n) \mathbf{x}_n$$

- Thus the MLE for  $\boldsymbol{\mu}_i$  is simply the mean of the points in class  $C_i$ .
- This result is easily generalized to the multiclass case. The MLE of the prior probability associated with class  $C_k$  is given by the fraction of the training set points  $N_k$ .

# Maximum Likelihood Solution

$$p(t | \pi, \mu_1, \mu_2, \Sigma) = \prod_{n=1}^N [\pi \mathcal{N}(x_n | \mu_1, \Sigma)]^{t_n} [(1-\pi) \mathcal{N}(x_n | \mu_2, \Sigma)]^{1-t_n}$$

- Consider finally the maximization with respect to  $\Sigma$ . The terms of interest in the log-likelihood are:

$$\begin{aligned} & -\frac{1}{2} \sum_{n=1}^N t_n \ln |\Sigma| - \frac{1}{2} \sum_{n=1}^N t_n (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_1) \\ & - \frac{1}{2} \sum_{n=1}^N (1-t_n) \ln |\Sigma| - \frac{1}{2} \sum_{n=1}^N (1-t_n) (\mathbf{x}_n - \boldsymbol{\mu}_2)^T \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_2) \\ & = -\frac{N}{2} \ln |\Sigma| - \frac{N}{2} \text{Tr}(\Sigma^{-1} S), \text{ where } S = \frac{N_1}{N} S_1 + \frac{N_2}{N} S_2 \end{aligned}$$

$$S_1 = \frac{1}{N_1} \sum_{n \in C_1} (\mathbf{x}_n - \boldsymbol{\mu}_1)(\mathbf{x}_n - \boldsymbol{\mu}_1)^T, S_2 = \frac{1}{N_2} \sum_{n \in C_2} (\mathbf{x}_n - \boldsymbol{\mu}_2)(\mathbf{x}_n - \boldsymbol{\mu}_2)^T$$

- Using the result for the MLE for a Gaussian distribution, we obtain  $\Sigma = S$ . This result is a weighted average of the covariance matrices associated with each class separately.

# Maximum Likelihood Solution

$$\pi = \frac{N_1}{N_1 + N_2}, \quad \boldsymbol{\mu}_1 = \frac{1}{N_1} \sum_{n=1}^N t_n \mathbf{x}_n, \quad \boldsymbol{\mu}_2 = \frac{1}{N_2} \sum_{n=1}^N (1-t_n) \mathbf{x}_n$$

$$\Sigma = \frac{N_1}{N} \mathbf{S}_1 + \frac{N_2}{N} \mathbf{S}_2, \quad \mathbf{S}_1 = \frac{1}{N_1} \sum_{n \in C_1} (\mathbf{x}_n - \boldsymbol{\mu}_1)(\mathbf{x}_n - \boldsymbol{\mu}_1)^T, \quad \mathbf{S}_2 = \frac{1}{N_2} \sum_{n \in C_2} (\mathbf{x}_n - \boldsymbol{\mu}_2)(\mathbf{x}_n - \boldsymbol{\mu}_2)^T$$

- These results are easily extended to the  $K$  class problem to obtain the corresponding MLE solutions for the parameters (for Gaussian class-conditionals with a shared covariance matrix).
- As we have seen before, MLE is not robust to outliers.
  - Fitting Gaussian distributions to the classes is not robust to outliers, because the MLE of a Gaussian is not robust.

# Discrete Features: Naive Bayes Assumption

- Here we consider the case of discrete feature values  $x_i$ . Consider binary feature values  $x_i \in \{0, 1\}, i = 1, \dots, D$ .
- For  $D$  inputs, a general distribution corresponds to a table of  $2^D$  numbers for each class, containing  $2^D - 1$  independent variables (due to the summation constraint).
- This grows exponentially with the number  $D$  of features.
- Consider the **Naive Bayes Assumption**: the feature values are treated as independent, conditioned on the class  $C_k$ .
- Thus we have class-conditional distributions of the form

$$p(\mathbf{x} | C_k) = \prod_{i=1}^D \mu_{ki}^{x_i} (1 - \mu_{ki})^{1-x_i}$$

# Discrete Features

---

$$p(\mathbf{x} | C_k) = \prod_{i=1}^D \mu_{ki}^{x_i} (1 - \mu_{ki})^{1-x_i}$$

- There are  $D$  independent parameters for each class.
- Here  $\mu_{ki}$  is the probability the  $i^{\text{th}}$  feature takes the value 1.
- Substituting into  $a_k = \ln(p(\mathbf{x} | C_k)p(C_k))$  gives:

$$a_k(\mathbf{x}) = \sum_{i=1}^D \left\{ x_i \ln \mu_{ki} + (1 - x_i) \ln (1 - \mu_{ki}) \right\} + \ln p(C_k)$$

- These are linear functions of the input values  $x_i$ .
  - For  $K = 2$ , we can alternatively consider the logistic sigmoid
- $$p(C_1 | \mathbf{x}) = \sigma(a), a(\mathbf{x}) = \ln \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_2)p(C_2)}$$
- Analogous results are obtained for discrete variables each of which can take  $M > 2$  states.

# Exponential Family

- We have seen that for both Gaussian and discrete inputs, the posterior class probabilities are given by generalized linear models with logistic sigmoid ( $K = 2$ ) or softmax ( $K \geq 2$  classes) activation functions.
- This result is general for the class-conditional densities  $p(x|C_k)$  that are members of the exponential family.
- For the exponential family, the distribution of  $x$  can be written in the form

$$p(x | \lambda_k) = h(x)g(\lambda_k)\exp\left\{\lambda_k^T u(x)\right\}$$

- Here, we consider the case  $u(x) = x$ .

# Exponential Family

$$p(\mathbf{x} | \boldsymbol{\lambda}_k) = h(\mathbf{x})g(\boldsymbol{\lambda}_k)\exp\left\{\boldsymbol{\lambda}_k^T \mathbf{u}(\mathbf{x})\right\}$$

- With  $\mathbf{u}(\mathbf{x}) = \mathbf{x}$  and by introducing a scaling parameter  $s$  as  $p(\mathbf{x} | s) = (1/s)f(\mathbf{x}/s)$ , we obtain the restricted set of exponential family class-conditional densities of the form:

$$p(\mathbf{x} | \boldsymbol{\lambda}_k, s) = \frac{1}{s} h\left(\frac{1}{s}\mathbf{x}\right)g(\boldsymbol{\lambda}_k)\exp\left\{\frac{1}{s}\boldsymbol{\lambda}_k^T \mathbf{x}\right\}$$

- Here each class has its own parameter  $\boldsymbol{\lambda}_k$  but all classes share the same scale parameter  $s$ .
- For the 2 class problem, we substitute this expression for the class-conditional densities into

$$a(\mathbf{x}) = \ln \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_2)p(C_2)}$$

# Exponential Family

---

$$p(\mathbf{x} | \boldsymbol{\lambda}_k, s) = \frac{1}{s} h\left(\frac{1}{s} \mathbf{x}\right) g(\boldsymbol{\lambda}_k) \exp\left\{\frac{1}{s} \boldsymbol{\lambda}_k^T \mathbf{x}\right\} \quad a(\mathbf{x}) = \ln \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_2)p(C_2)}$$

- The posterior class probability is again given by a logistic sigmoid acting on a linear function  $a(\mathbf{x})$

$$p(C_1 | \mathbf{x}) = \sigma(a(\mathbf{x}))$$

which is given by

$$a(\mathbf{x}) = \frac{1}{s} (\boldsymbol{\lambda}_1 - \boldsymbol{\lambda}_2)^T \mathbf{x} + \ln g(\boldsymbol{\lambda}_1) - \ln g(\boldsymbol{\lambda}_2) + \ln p(C_1) - \ln p(C_2)$$

- For the  $K$  –class problem, we substitute the class-conditionals

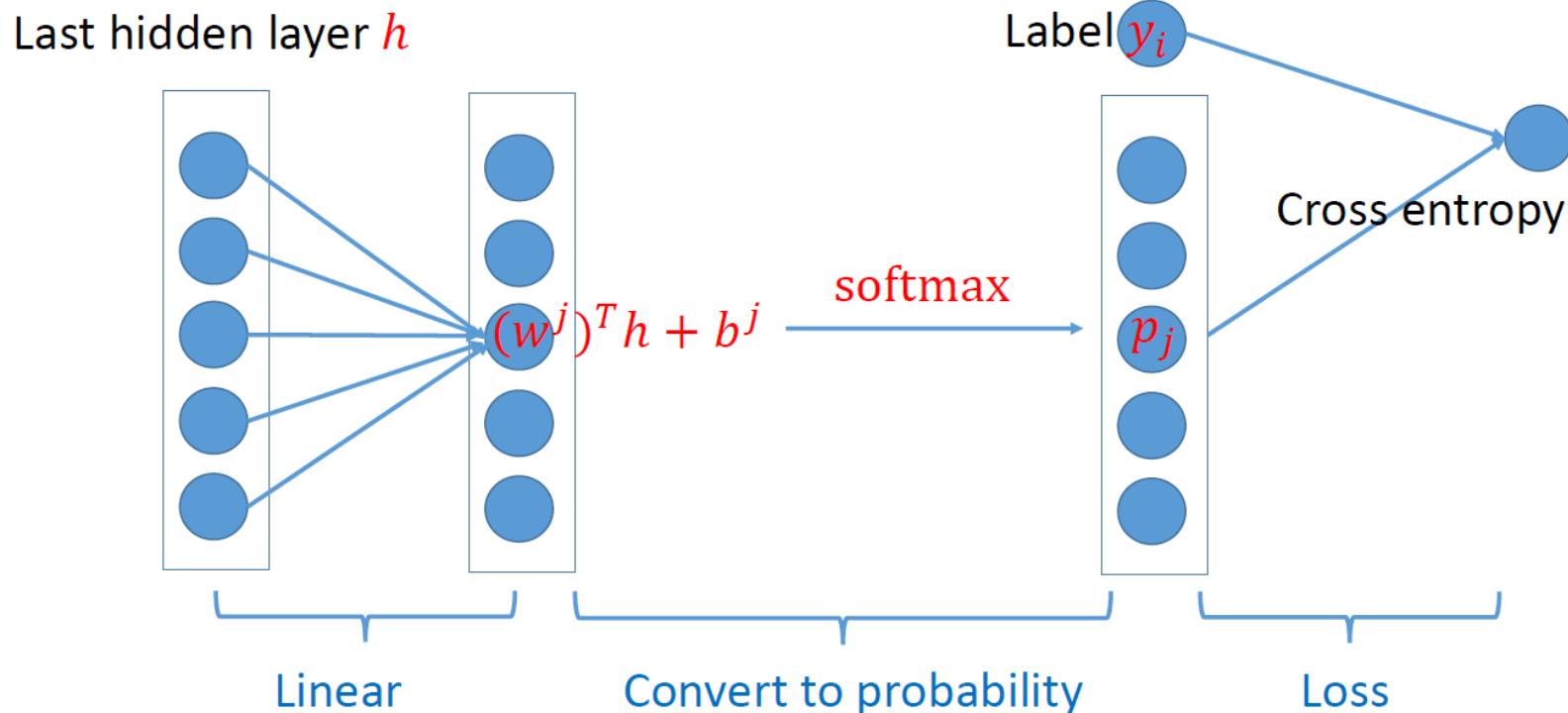
directly in  $a_k(\mathbf{x}) = \ln(p(\mathbf{x} | C_k)p(C_k))$  and  $p(C_k | \mathbf{x}) = \frac{e^{a_k}}{\sum_j e^{a_j}}$  to obtain:

$$a_k(\mathbf{x}) = \frac{1}{s} \boldsymbol{\lambda}_k^T \mathbf{x} + \ln g(\boldsymbol{\lambda}_k) + \ln p(C_k)$$

which is linear in  $\mathbf{x}$ .

# Cross Entropy for Conditional Distribution

- Let  $p_{data}(y|x)$  denote the empirical distribution of the data.
- The negative log-likelihood:  $-\frac{1}{N} \sum_{i=1}^N \log p(y = y_i | x_i) = -\mathbb{E}_{p_{data}} \log p(y|x)$  is the cross entropy between  $p_{data}$  and the model output  $p$ .
- Write this as KL minimization:  $D(p_{data} || p) = \mathbb{E}_{p_{data}} [\log p_{data}] - \mathbb{E}_{p_{data}} [\log p]$ . The 1<sup>st</sup> term (entropy) is const. & the 2<sup>nd</sup> the cross entropy.



---

# **Stochastic Gradient Descent Robbins-Monro Algorithm**

# Stochastic Gradient Descent

---

- **Stochastic gradient descent** is a gradient descent optimization method for minimizing an objective function written as a sum of differentiable functions.
- Consider minimizing an objective function that has the form of a sum ( $M$  –estimators)

$$f(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \underbrace{f(\boldsymbol{\theta}, z_i)}_{\text{Loss Function}}, z_i = (\mathbf{x}_i, y_i)$$

- The parameter  $\theta$  is to be estimated and each summand function is related with the  $i$  –th observation in the data set (used for training).

- Typical problems are the maximization (stationary points) of the likelihood function,

$$f(\boldsymbol{\theta}, z_i) = -\log p(y_i | \mathbf{x}_i, \boldsymbol{\theta})$$

- Minimization of a loss function, etc.

$$f(\boldsymbol{\theta}, \mathbf{z}_i) = L(y_i, \mathbf{h}(\mathbf{x}_i, \boldsymbol{\theta})) = L(y_i, \hat{y}_i), \hat{y}_i \equiv h(\mathbf{x}_i, \boldsymbol{\theta})$$

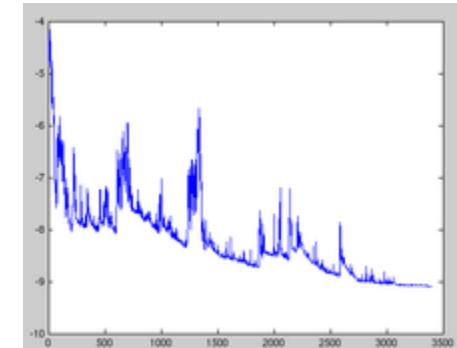
# Stochastic Gradient Descent

- A batch gradient descent method performs the following iterations

$$\theta_{k+1} = \theta_k - \eta \sum_{i=1}^N \nabla f(\theta, z_i)$$

- $\eta$  is the *learning rate*. Fluctuations in the objective function with gradient steps are shown here
- Evaluating the sums of gradients is expensive.
- Stochastic gradient descent samples a subset of summand functions at every step. This is very effective in the case of large-scale machine learning problems.
- For example passing through each data point in the training set:

$$\theta_{k+1} = \theta_k - \eta \nabla f(\theta, z_i)$$



- Zinkevich, M. (2003). [Online convex programming and generalized infinitesimal gradient ascent](#). In *Intl. Conf. on Machine Learning*, pp. 928-936.
- Cesa-Bianchi, N. and G. Lugosi (2006). [Prediction, learning, and games](#). Cambridge University Press.
- [Nemirovski, A.](#) and D. Yudin (1978). On Cezari's convergence of the [steepest descent method](#) for approximating saddle points of convex-concave functions. *Soviet Math. Dokl.* 19.
- Sra et al. [Optimization for Machine Learning](#), 2012.

# Stochastic Gradient Descent

---

- One can use a running average of the parameter

$$\bar{\theta}_k = \frac{1}{k} \sum_{i=1}^k \theta_i$$

- Or use the Polyak-Ruppert averaging

$$\bar{\theta}_k = \bar{\theta}_{k-1} - \frac{1}{k} (\theta_{k-1} - \theta_k)$$

- Spall, J. (2003). *Introduction to Stochastic Search and Optimization: Estimation, Simulation, and Control*. Wiley.
- Kushner, H. and G. Yin (2003). *Stochastic approximation and recursive algorithms and applications*. Springer.

# Stochastic Optimization

---

- Stochastic optimization algorithms attempt to find zeroes or extrema of functions which cannot be computed directly, but only estimated via noisy observations. Mathematically, this refers to solving:

$$\min_{\theta} f(\theta) = \mathbb{E}[f(\theta, z)]$$

- The objective is to find the parameter  $\theta$ , which minimizes  $f(\theta)$  for some unknown random variable,  $z$ .
- $f(\theta)$  cannot be computed exactly, but instead approximated via simulation. This situation is modeled by the function  $f(\theta, z)$ , where  $z$  represents the noise and is a random variable.
- The objective is then to minimize the expectation of  $f(\theta, z)$ . Our earlier objective function can be considered as a sampled approximation of the expectation.
- The first algorithm of this kind is the Robbins-Monro algorithm.

- [Robbins H., Monro, S. \(1951\). A Stochastic Approximation Method. \*The Annals of Mathematical Statistics\* 22 \(3\): 400.](#)

# **Robbins Monro Algorithm**

---

- Robbins–Monro algorithm first introduced to solve a root finding problem where the function is represented as expected value.
- Assume that we have a function  $M(\theta)$ , and a constant  $\alpha$ , such that the equation  $M(\theta) = \alpha$  has a unique root at  $\theta = \theta^*$ .
- It is assumed that while we cannot directly observe the function  $M(\theta)$ , we can instead obtain measurements of the random variable  $N(\theta)$  where  $\mathbb{E}[N(\theta)] = M(\theta)$ . The structure of the algorithm is to then generate iterates of the form:

$$\theta_{n+1} = \theta_n - \eta_n (N(\theta_n) - \alpha), \quad \eta_n \text{ sequence of } > 0 \text{ stepsizes}$$

- The sequence  $\theta_n$  converges in  $L_2$  to  $\theta^*$ , provided that:
  - $N(\theta)$  is uniformly bounded
  - $M(\theta)$  is nondecreasing
  - $M'(\theta^*)$  exists and is positive
  - The following hold:  $\sum_{n=0}^{\infty} \eta_n = \infty, \sum_{n=0}^{\infty} \eta_n^2 < \infty$
- A particular choice is:  $\eta_n = a/n$ , for some  $a > 0$ .

# Robbins Monro Algorithm

---

- Another choice of the learning rate schedule is:

$$\eta_k = \frac{1}{(\tau_0 + k)^\kappa}, \tau_0 \geq 0, \kappa \in (0.5, 1]$$

- The need to adjust these tuning parameters is one of the main drawback of stochastic optimization.
- A simple approach is as follows:
  - store an initial subset of the data, and try a range of  $\eta$  values on this subset;
  - then choose the one that results in the fastest decrease in the objective and apply it to all the rest of the data.
  - Convergence is not certain but the algorithm can be terminated when the performance improvement on a hold-out set plateaus (early stopping).

- [Bottou, L. \(1998\). Online algorithms and stochastic approximations](#). In D. Saad (Ed.), *Online Learning and Neural Networks*. Cambridge ([presentations](#))
- Bach, F. and E. Moulines (2011). [Nonasymptotic analysis of stochastic approximation algorithms for machine learning](#). In *NIPS*.
- [Bottou, L.](#) (2007). [Learning with large datasets](#) (nips tutorial).

# Per-Parameter Step Sizes

---

- One drawback of SGD is that it uses the same step size for all parameters. This can be extended using a method known as [adagrad \(adaptive gradient\)](#).
- In particular, if  $\theta_i(k)$  is parameter  $\theta_i$  at time  $k$ , and  $g_i(k)$  is its gradient, then we make an update as follows:

$$\theta_i(k+1) = \theta_i(k) - \eta \frac{g_i(k)}{\tau_0 + \sqrt{s_i(k)}}, \quad s_i(k) = s_i(k-1) + g_i(k)^2$$

- Since  $s_i(k)$  is the  [\$\ell\_2\$  norm](#) of previous derivatives, extreme parameter updates get damped, while parameters that get few or small updates receive higher learning rates.
- The result is *a per-parameter step size that adapts to the curvature of the loss function.*

- [Duchi, J., E. Hazan, and Y. Singer](#) (2010). [Adaptive Subgradient Methods for Online Learning and Stochastic Optimization](#). In *Proc. of the Workshop on Computational Learning Theory*
- [Schaul, T., S. Zhang, and Y. LeCun](#) (2012). [No more pesky learning rates](#). Technical report, Courant Institute of Mathematical Sciences.

# The LMS Algorithm For Regression

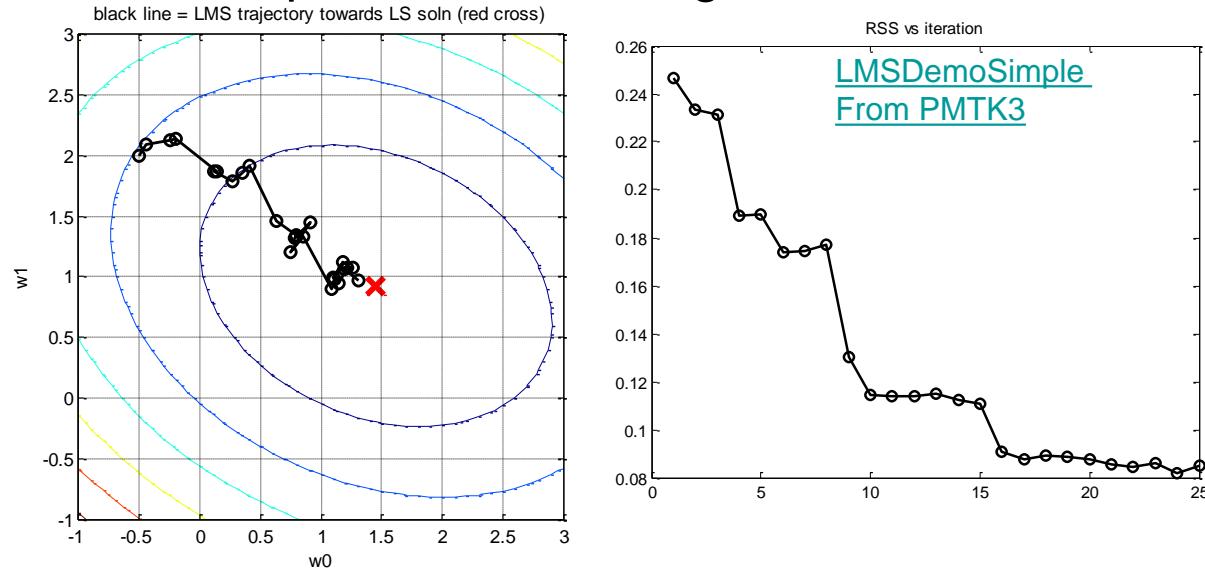
- Let us recall the MLE for linear regression in an online fashion (see earlier lecture on Regression). The online gradient at iteration  $k$  is

$$\mathbf{g}_k = \mathbf{x}_i (\boldsymbol{\theta}_k^T \mathbf{x}_i - y_i)$$

- The gradient acts like an error signal. After computing the gradient, we take a step along it as follows:

$$\boldsymbol{\theta}_{(k+1)} = \boldsymbol{\theta}_{(k)} - \eta_k (\hat{y}_k - y_k) \mathbf{x}_k, \hat{y}_k = \boldsymbol{\theta}_k^T \mathbf{x}_k$$

- This algorithm is called the **least mean squares (LMS)** and is an example of the SGD algorithm.



- LMS may require multiple passes through the data to find the optimum.

- The recursive least squares algorithm, based on the Kalman filter and which uses 2<sup>nd</sup>-order information, finds the optimum in a single pass (see Murphy, Section 18.2.3)

---

# The Perceptron Algorithm

# The Rosenblatt Perceptron Algorithm

---

- This is a two-class model in which  $x$  is mapped to a feature  $\phi(x)$ , that is used to construct a generalized linear model

$$y(x) = f(\mathbf{w}^T \phi(x))$$

where the activation function  $f$  is:

$$f(a) = \begin{cases} +1 & a \geq 0 \\ -1 & a < 0 \end{cases}$$

- The vector  $\phi(x)$  includes a bias component  $\phi_0(x) = 1$ .
- We use target  $t = +1$  for class  $C_1$  and  $t = -1$  for class  $C_2$ .
- When using as error function the total number of misclassified patterns, the error is a piecewise constant function of  $\mathbf{w}$ , with discontinuities wherever a change in  $\mathbf{w}$  causes the decision boundary to move across one of the data points.
  - Gradient optimization methods cannot be applied, because the gradient is zero almost everywhere.
  - Rosenblatt, F. (1962). [\*Principles of Neurodynamics: Perceptrons and the Theory of Brain Mechanisms\*](#). Spartan.

# The Rosenblatt Perceptron Algorithm

---

- We consider an alternative error function (perceptron criterion).
- We are seeking a weight vector  $w$  such that  $x_n$  in class  $C_1$  will have  $w^T \phi(x_n) > 0$ , whereas  $x_n$  in  $C_2$  have  $w^T \phi(x_n) < 0$ .
- Using the  $t \in \{1, -1\}$  target coding scheme it follows that we would like all patterns to satisfy  $w^T \phi(x_n) t_n > 0$ .
- *The perceptron associates zero error with a correctly classified pattern, whereas for a misclassified  $x_n$  it minimizes  $-w^T \phi(x_n) t_n$ .*
- The perceptron criterion is

$$E_P(w) = - \sum_{n \in \mathcal{M}} w^T \phi_n t_n$$

$\mathcal{M}$  is the set of all misclassified patterns.

# The Rosenblatt Perceptron Algorithm

---

- The contribution to the error from a misclassified pattern is a linear in  $\mathbf{w}$  in regions of  $\mathbf{w}$ -space where the pattern is misclassified and zero in regions where it is correctly classified.
- The total error function is piecewise linear

$$E_P(\mathbf{w}) = - \sum_{n \in \mathcal{M}} \mathbf{w}^T \phi_n t_n$$

- We apply the stochastic gradient descent algorithm:

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_P(\mathbf{w}) = \mathbf{w}^{(\tau)} + \eta \phi_n t_n$$

where  $\eta$  is the learning rate parameter.

- Since the perceptron function  $y(\mathbf{x}, \mathbf{w})$  is unchanged if we multiply  $\mathbf{w}$  by a const., we set  $\eta = 1$  without loss of generality.
  - During training, the set of patterns that are misclassified will change.

# The Rosenblatt Perceptron Algorithm

---

- The perceptron learning algorithm has a simple interpretation

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} + \phi_n t_n$$

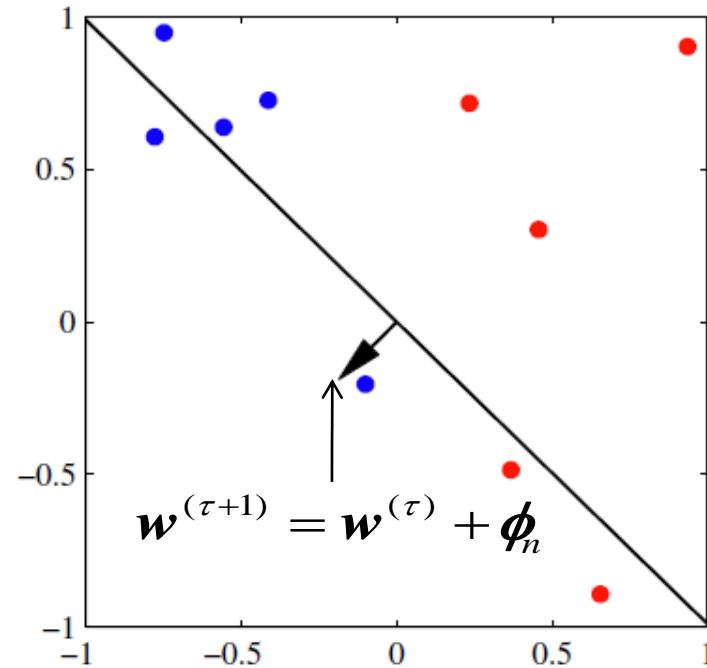
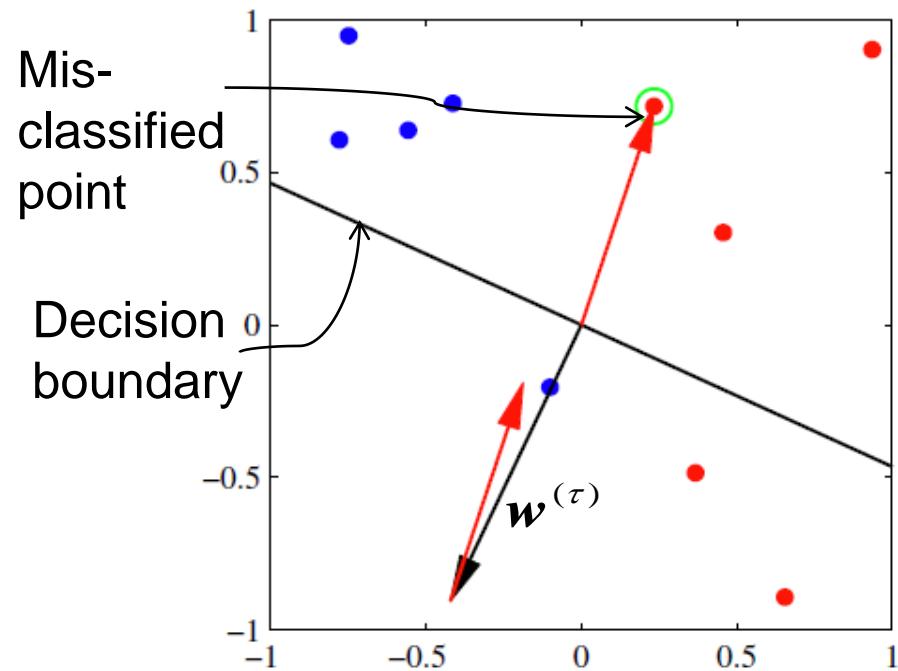
- We cycle through the training patterns, and for each  $x_n$  we evaluate the perceptron function  $y(x) = f(\mathbf{w}^T \phi(x))$ .

$$f(a) = \begin{cases} +1 & a \geq 0 \\ -1 & a < 0 \end{cases}$$

- If the pattern is correctly classified, then the weight vector remains unchanged.
- If the pattern is incorrectly classified, then for class  $C_1$  we add  $\phi(x_n)$  onto the current estimate of weight vector  $\mathbf{w}$  while for  $C_2$  we subtract  $\phi(x_n)$  from  $\mathbf{w}$ .

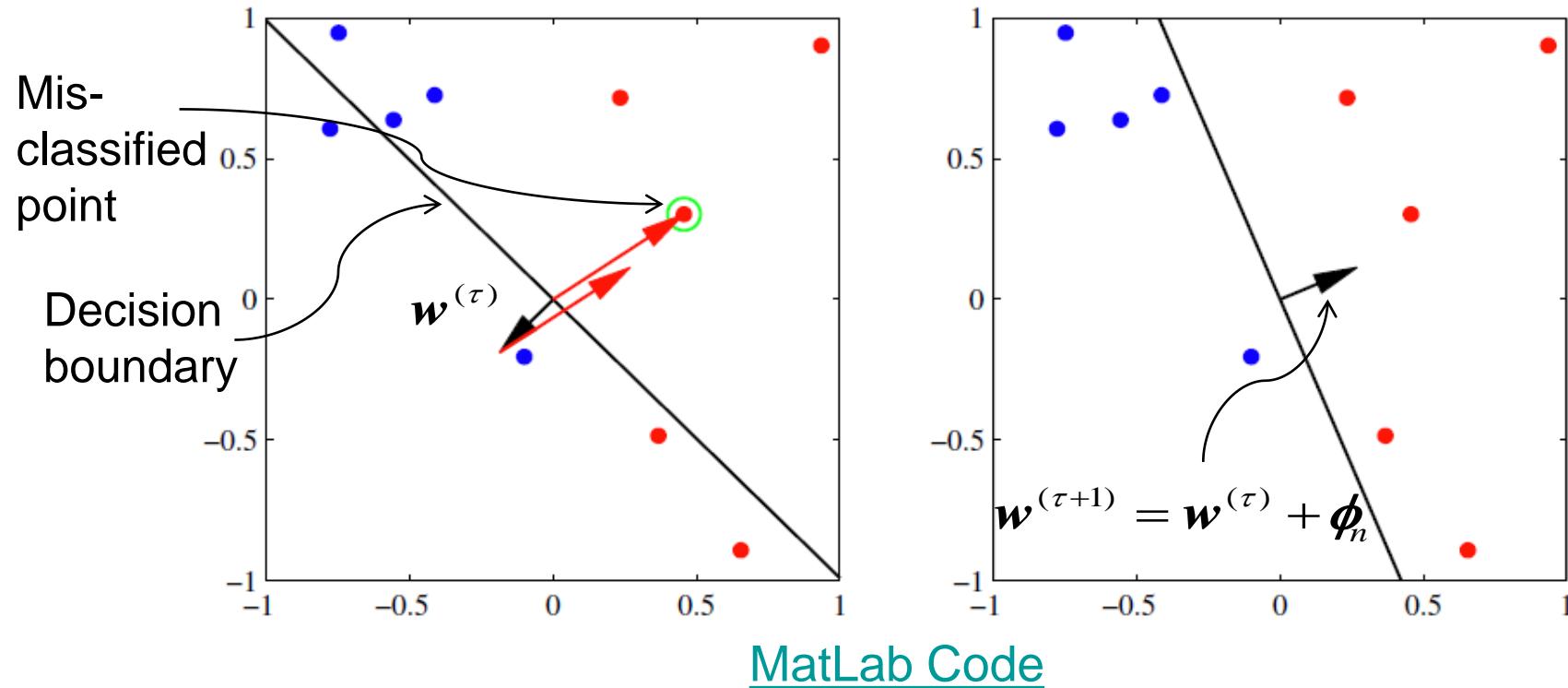
# The Rosenblatt Perceptron Algorithm

- Data points from two classes (red and blue) in a feature space ( $\phi_1, \phi_2$ ).
- The parameter vector  $w^{(\tau)}$  points towards the red class.
- The data point shown is misclassified and so its feature vector is added to the current weight vector, giving the new decision boundary on the right.



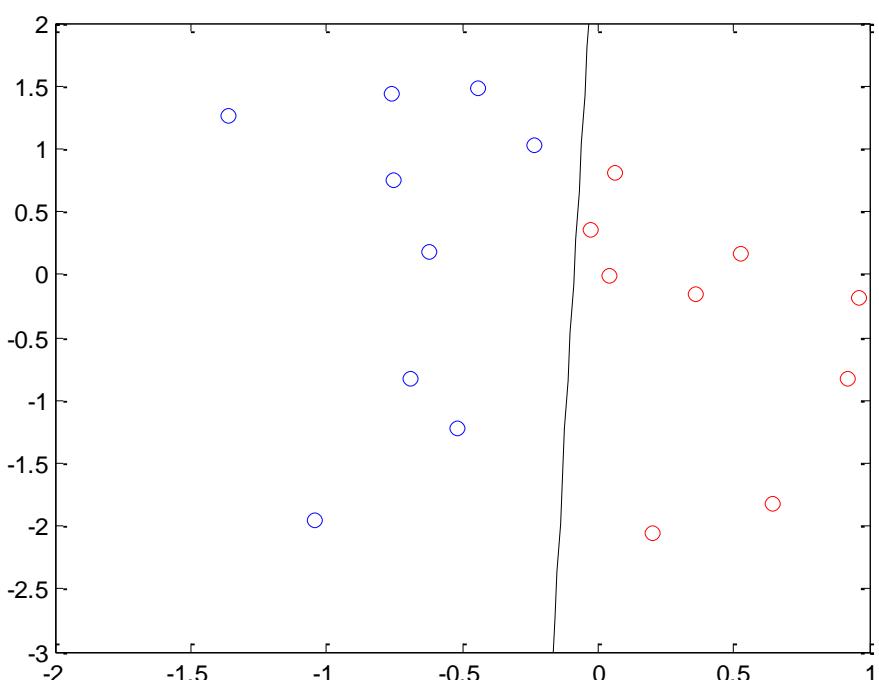
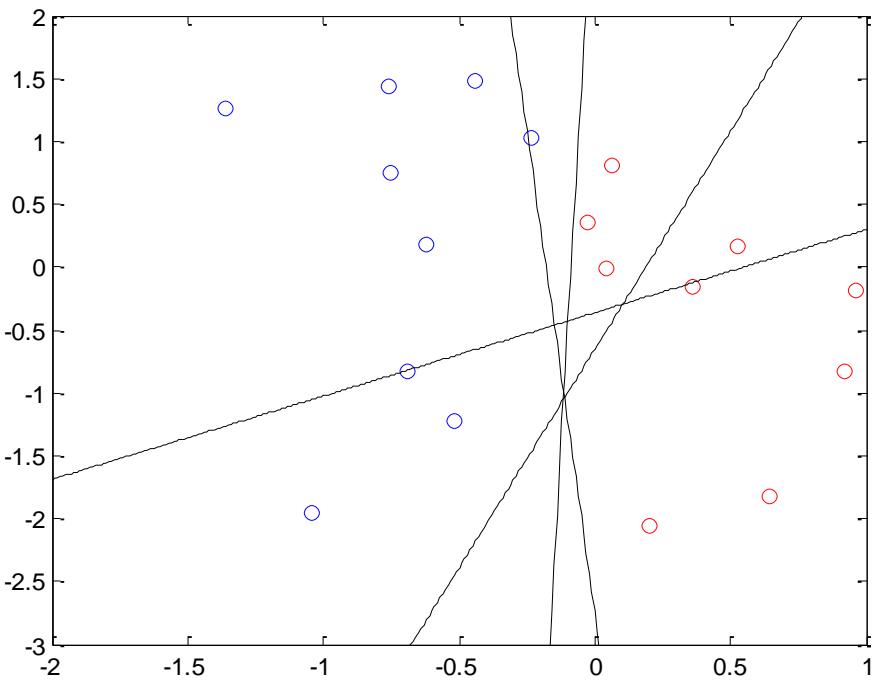
# The Rosenblatt Perceptron Algorithm

- The left plot shows the next misclassified point.
- Its feature vector is added to the weight vector giving the decision boundary shown in the right plot for which all data points are correctly classified.



# The Rosenblatt Perceptron Algorithm

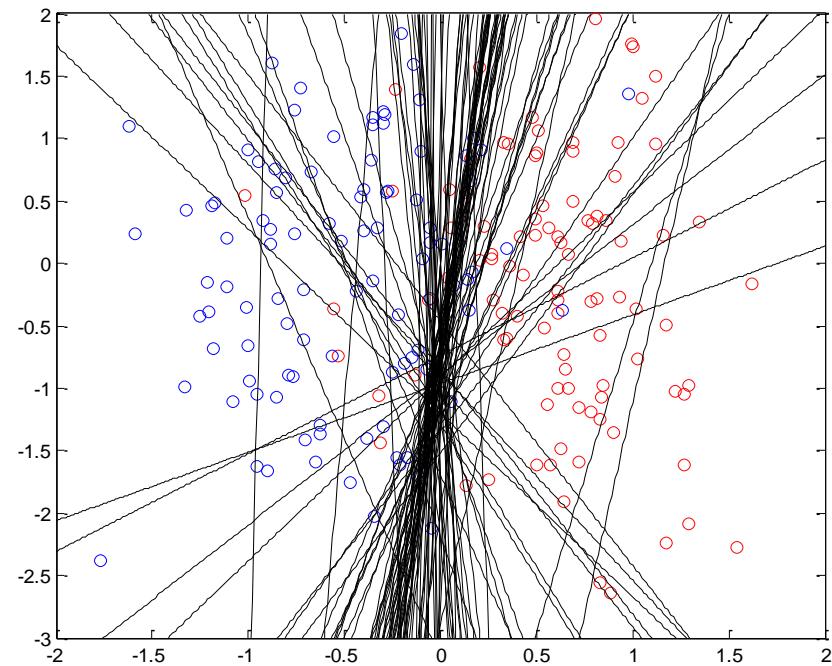
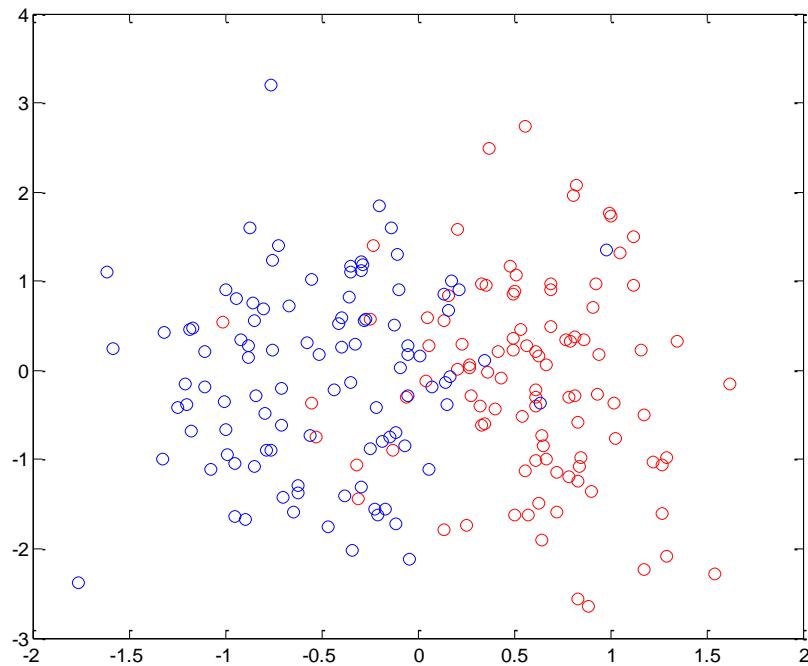
- Demonstration of the algorithm using the data set *data\_RP\_10* (very fast convergence)



MatLab Code

# The Rosenblatt Perceptron Algorithm

- ❑ Use data set *data\_RP\_100* (very slow convergence)
- ❑ Weakness: the algorithm will never converge if the two classes **overlap**, which means we cannot find a decision boundary that completely separates them.



MatLab Code

# The Rosenblatt Perceptron Algorithm

---

- In each update in the perceptron learning algorithm, the contribution to the error from a misclassified pattern will be reduced because from  $\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} + \eta \phi_n t_n$  we have

$$-\mathbf{w}^{(\tau+1)T} \phi_n t_n = -\mathbf{w}^{(\tau)T} \phi_n t_n - (\phi_n t_n)^T \phi_n t_n < -\mathbf{w}^{(\tau)T} \phi_n t_n$$

where we set  $\eta = 1$ .

- This does not imply that the contribution to the error from the other misclassified patterns will have been reduced.
- Also the change in  $w$  may have caused some previously correctly classified patterns to become misclassified.

The perceptron learning rule is not guaranteed to reduce the total error function at each stage

# The Perceptron Convergence Theorem

---

- For linearly separable training data set, the perceptron learning algorithm is guaranteed to find an exact solution in a finite number of steps.
- The number of steps to achieve convergence could be substantial.
- Until convergence is achieved, one cannot distinguish a nonseparable problem and one that is slow to converge.
- There may be many solutions: initialization of the parameters and the order of going through the data points affects the solution obtained.

For data sets that are not linearly separable, the perceptron learning algorithm will never converge

# The Perceptron Convergence Theorem

---

- The perceptron does not provide probabilistic outputs.
- It cannot be applied to  $K > 2$  classes.
- Like with all other classification algorithms examined up to now, it is based on linear combinations of fixed basis functions.

- Rosenblatt, F. (1962). *Principles of Neurodynamics: Perceptrons and the Theory of Brain Mechanisms*. Spartan.
- Block, H. D. (1962). *The perceptron: a model for brain functioning*. *Reviews of Modern Physics* **34**(1), 123–135. Reprinted in Anderson and Rosenfeld (1988).
- Nilsson, N. J. (1965). *Learning Machines*. McGraw- Hill. Reprinted as *The Mathematical Foundations of Learning Machines*, Morgan Kaufmann, (1990).
- Minsky, M. L. and S. A. Papert (1969). *Perceptrons*. MIT Press. Expanded edition 1990.
- Hertz, J., A. Krogh, and R. G. Palmer (1991). *Introduction to the Theory of Neural Computation*. Addison Wesley.
- Bishop, C. M. (1995a). *Neural Networks for Pattern Recognition*. Oxford University Press.
- Widrow, B. and M. E. Hoff (1960). *Adaptive switching circuits*. In *IRE WESCON Convention Record*, Volume 4, pp. 96–104. Reprinted in Anderson and Rosenfeld (1988).
- Widrow, B. and M. A. Lehr (1990). *30 years of adaptive neural networks: perceptron, madeline, and backpropagation*. *Proceedings of the IEEE* **78**(9), 1415–1442.

# A Bayesian Perspective

---

- Another approach to online learning is to adopt a Bayesian view. We just apply Bayes rule recursively:

$$p(\boldsymbol{\theta}|\mathcal{D}_{1:k}) \propto p(\mathcal{D}_k|\boldsymbol{\theta})p(\boldsymbol{\theta}|\mathcal{D}_{1:k-1})$$

- This has the advantage of returning a posterior instead of just a point estimate. It also allows for the online adaptation of hyper-parameters.
- Finally, it is quicker than SGD. Note that *by modeling the posterior variance of each parameter, we effectively associate a different learning rate for each parameter, which is a simple way to model the curvature of the space.*
- These variances can then be adapted using the usual rules of probability theory.

▪ [de Freitas, N., M. Niranjan, and A. Gee](#) (2000). [Hierarchical Bayesian models for regularisation in sequential learning](#). *Neural Computation* 12(4), 955–993.

# A Bayesian Perspective

---

- Getting 2<sup>nd</sup>-order optimization methods to work online is tricky.<sup>1-4</sup>
- One can use the Kalman filter to fit a linear regression model online. Unlike the LMS algorithm, this converges to the optimal (offline) answer in a single pass over the data.
- An extension which can learn a robust non-linear regression model in an online fashion is described in [5]. For the GLM case, we can use an assumed density filter, where we approximate the posterior by a Gaussian with a diagonal covariance; the variance terms serve as a per-parameter step-size.
- Another approach is to use particle filtering [6] for sequentially learning a kernelized linear/logistic regression model.

1. [Schraudolph, N. N., J. Yu, and S. Gunter](#) (2007). [A Stochastic Quasi- Newton Method for Online Convex Optimization](#). In *AI/Statistics*, pp. 436–443.
2. [Sunehag, P., J. Trumpf, S. V. N. Vishwanathan, and N. N. Schraudolph](#) (2009). [Variable Metric Stochastic Approximation Theory](#). In *AI/Statistics*, pp. 560–566.
3. [Bordes, A., L. Bottou, and P. Gallinari](#) (2009, July). [Sgd-qn: Careful quasi-newton stochastic gradient descent](#). *J. of Machine Learning Research* 10, 1737–1754.
4. [Bordes, A., L. Bottou, P. Gallinari, J. Chang, and S. A. Smith](#) (2010). [Erratum: SGDQN is Less Careful than Expected](#). *J. of Machine Learning Research* 11, 2229–2240.
5. [Ting, J., A. D'Souza, S. Vijayakumar, and S. Schaal](#) (2010). [Efficient learning and feature selection in high-dimensional regression](#). *Neural Computation* 22(4), 831–886.
6. Andrieu, C., N. de Freitas, and A. Doucet (2000). [Sequential Bayesian estimation and model selection for dynamic kernel machines](#). Technical report, Cambridge Univ.

---

# *Probabilistic Discriminative Models*

# Probabilistic Discriminative Models

---

- For the two-class problem, the posterior probability of class  $C_k$  can be written as a logistic sigmoid acting on a linear function of  $x$ , for a wide choice of  $p(x|C_k)$ .
- For the multiclass problem, the posterior probability of  $C_k$  is given by a softmax transformation of a linear function of  $x$ .
- For specific choices of  $p(x|C_k)$ , we used MLE to determine the parameters of the densities and the class priors  $p(C_k)$  and then used Bayes' theorem to find the posterior class probabilities.
- An alternative approach that we follow here is to *use the functional form of the generalized linear model explicitly and determine its parameters directly by using MLE*.
- There is an efficient algorithm finding such solutions known as iterative reweighted least squares, or IRLS.

# Probabilistic Discriminative Models

---

- Indirect approach – Generative Modeling: find the parameters of the generalized linear model, by fitting class-conditional densities and class priors separately and then applying Bayes' theorem.
  - With this model, we *can generate synthetic data by drawing values of  $x$  from the marginal  $p(x)$ .*
- Direct approach – Discriminative Modeling: we are maximizing a likelihood function defined through  $p(C_k|x)$ , which represents a form of discriminative training.
- An advantage of the discriminative approach is that there are typically fewer adaptive parameters to be determined.
- It leads to improved predictive performance when the class-conditional density assumptions are a poor approximation to the true distributions.

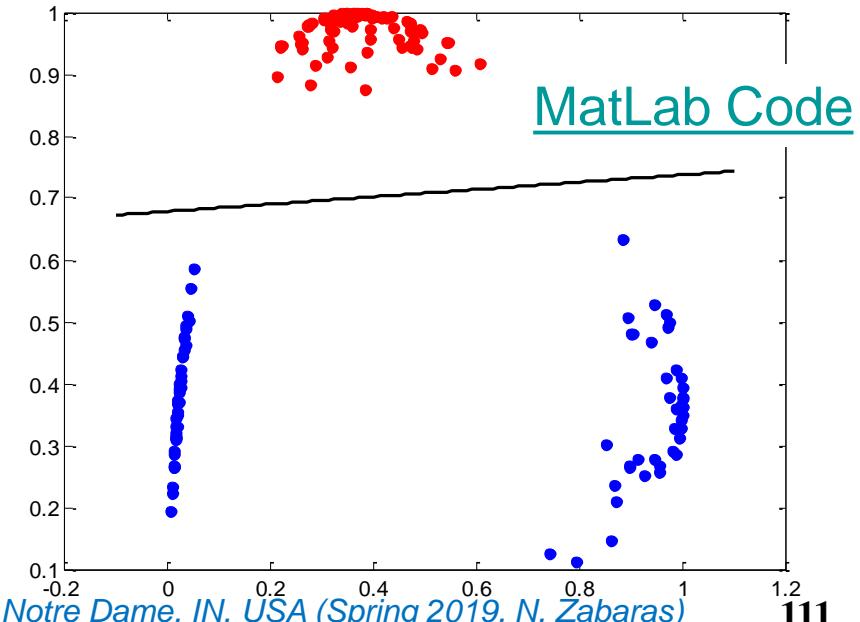
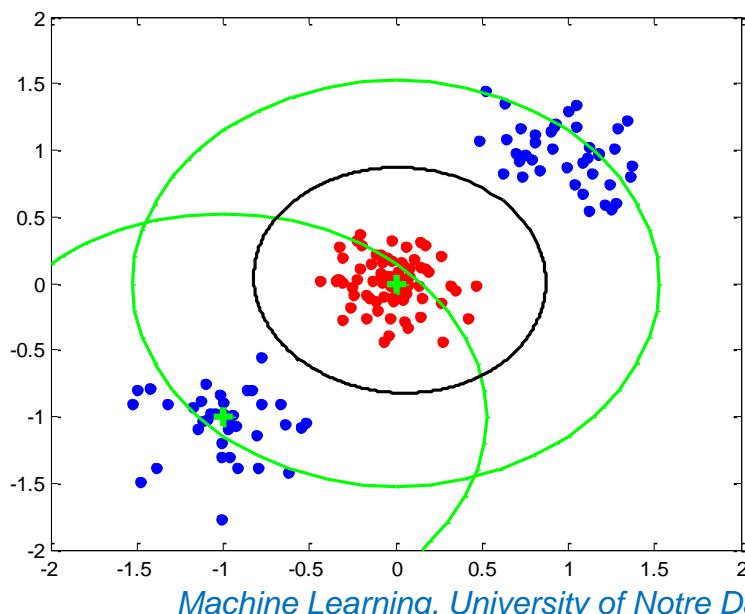
# Fixed Basis Functions

---

- Before we start with the logistic regression model, let us first make a fixed nonlinear transformation of the inputs using a vector of basis functions  $\phi(x)$ , then apply the algorithms discussed up to now.
- The resulting decision boundaries will be linear in the feature space  $\phi$ , and these correspond to nonlinear decision boundaries in the original  $x$  space.
- Classes that are linearly separable in the feature space  $\phi(x)$  need not be linearly separable in the original input space  $x$ .

# Nonlinear Basis in Classification Models

- Left: original space  $(x_1, x_2)$  with data from red & blue classes
- Two ‘Gaussian’ basis functions  $\phi_1(x), \phi_2(x)$  are defined in  $(x_1, x_2)$  with centers shown by the green crosses & contours by the green circles.
- Right: the corresponding feature space  $(\phi_1, \phi_2)$  together with the linear decision boundary obtained by a logistic regression model. This corresponds to a nonlinear decision boundary in  $(x_1, x_2)$  (black curve on the left).



# Fixed Basis Functions

---

- One of the basis functions is typically set to a constant, say  $\phi_0(x) = 1$ , so that  $w_0$  plays the role of a bias.
- For problems with significant overlap between the class-conditional densities  $p(x|C_k)$ , the posterior probabilities  $p(C_k|x)$  at least for some values of  $x$  may not be 0 or 1.
  - In such cases, the optimal solution is obtained by modeling the posterior probabilities accurately and then applying standard decision theory.
- The transformation  $\phi(x)$  cannot remove class overlap, it can increase the level of overlap, or create overlap where none existed in the  $x$  space.
- However, suitable choices of nonlinearity can make the process of modeling the posterior probabilities easier.

# **Fixed Basis Functions**

---

- Fixed basis function models have serious limitations.
- We should allow the basis functions to adapt to the data.
- Still models with fixed nonlinear basis functions are important as they introduce many of the key concepts needed for an understanding models with data-adaptive basis.

# Logistic Regression

- Consider a simple generalized linear problem with a two-class classification.
- In generative models we saw that the posterior probability of class  $C_1$  can be written as a logistic sigmoid acting on a linear function of the feature vector  $\phi$  so that

$$p(C_1|\phi) = y(\phi) = \sigma(w^T\phi)$$

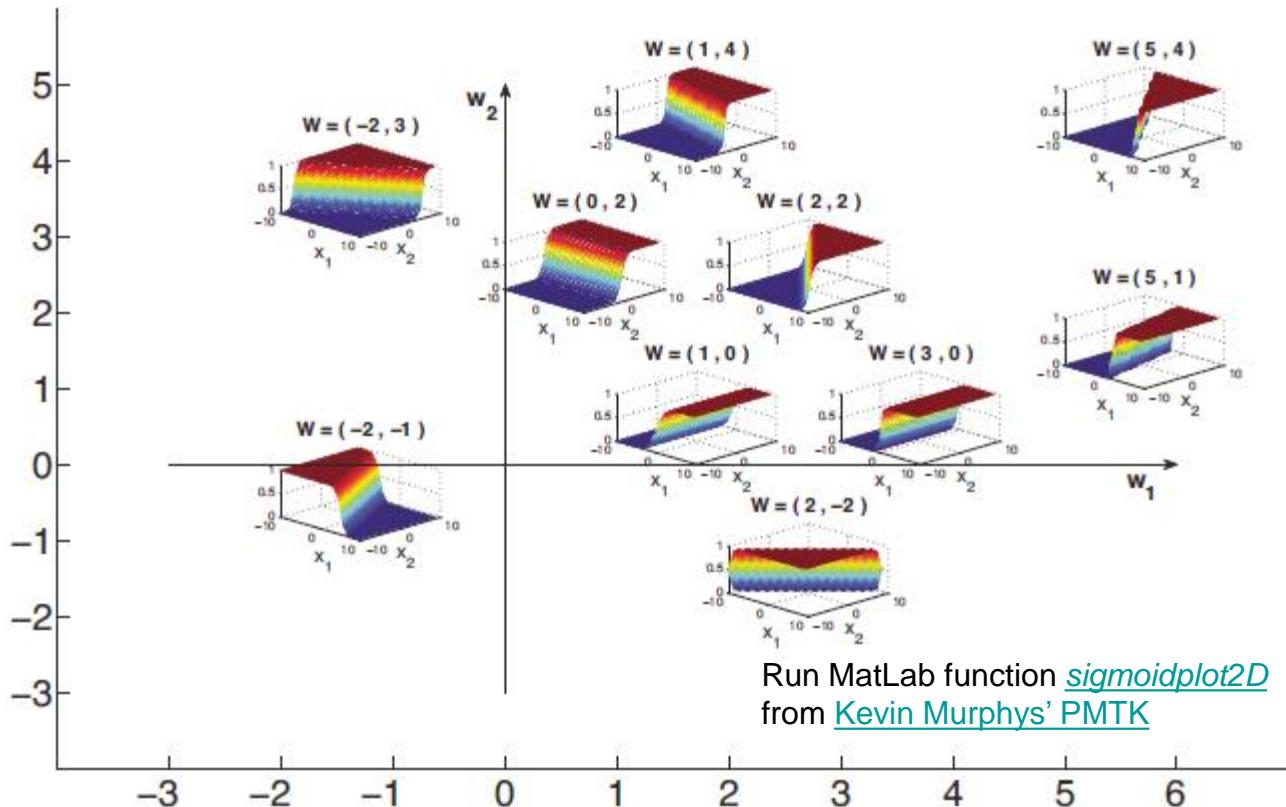
with  $p(C_2|\phi) = 1 - p(C_1|\phi)$ . Here  $\sigma(\cdot)$  is the logistic sigmoid function defined by

$$\sigma(a) = \frac{1}{1 + e^{-a}}$$

- This model is known as logistic regression – but *note its a model for classification not regression.*

# Plots of $p(C_1|\phi) = y(\phi) = \sigma(w^T\phi)$

- Plots of  $p(C_1|f) = \sigma(w_1x_1 + w_2x_2)$  for a 2D input and different  $w$ .
- If we threshold these probabilities at 0.5, we induce a linear decision boundary whose normal is  $w$ .
- To the right of this have  $\sigma(w^T x) > 0$  and to the left  $\sigma(w^T x) < 0$ .



# Logistic Regression

$$p(C_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi)$$

- For an  $M$  –dimensional feature space  $\phi$ , this model has  $M$  adjustable parameters.
- If we had fitted Gaussian class conditional densities using MLE, we would have used
  - $2M$  parameters for the two means and
  - $M(M + 1)/2$  parameters for the (shared) covariance.
  - Together with the class prior  $p(C_1)$ , this gives a total of  $M(M + 5)/2 + 1$  parameters, i.e.
  - $\mathcal{O}(M^2)$ , in contrast to  $\mathcal{O}(M)$  dependence on  $M$  in logistic regression.
- For large  $M$ , there is advantage in working with logistic regression.
- We now use MLE to determine the parameters in logistic regression. We will use:

$$\frac{d\sigma}{da} = \sigma(1 - \sigma)$$

# Logistic Regression - Cross-Entropy Error

- For a data set  $\{\phi_n, t_n\}$ , where  $t_n \in \{0, 1\}$  and  $\phi_n = \phi(x_n)$ , with  $n = 1, \dots, N$ , the likelihood function can be written

$$p(\mathbf{t} | \mathbf{w}) = \prod_{n=1}^N y_n^{t_n} (1 - y_n)^{1-t_n}, \mathbf{t} = (t_1, \dots, t_N)^T, y_n = p(C_1 | \phi_n)$$

- Using this, we define the cross-entropy error function as:

$$E(\mathbf{w}) = -\ln p(\mathbf{t} | \mathbf{w}) = -\sum_{n=1}^N \left\{ t_n \ln y_n + (1 - t_n) \ln (1 - y_n) \right\}, y_n = \sigma(a_n), a_n = \mathbf{w}^T \phi_n$$

- Taking the gradient wrt  $\mathbf{w}$  and using  $\frac{d\sigma}{da} = \sigma(1 - \sigma)$  gives:

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n, y_n = \sigma(a_n), a_n = \mathbf{w}^T \phi_n$$

- The contribution to  $\nabla E(\mathbf{w})$  from data point  $n$  is given by the 'error'  $y_n - t_n$  between the target value and the prediction of the model, times the basis  $\phi_n$ . This is identical to the gradient of the sum-of-squares error function for linear regression.

# Logistic Regression - Sequential Update

- We use this result

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n, \quad y_n = \sigma(a_n), \quad a_n = \mathbf{w}^T \phi_n$$

for sequential algorithm in which data are presented one at a time.

- The weight vector is updated using

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n = \mathbf{w}^{(\tau)} - \eta (y_n - t_n) \phi_n$$

- Note that for linearly separable data (the hyperplane  $\sigma = 0.5$ ,  $\mathbf{w}^T \phi = 0$  separates the two classes), any decision boundary separating the two classes needs to satisfy:

$$\mathbf{w}^T \phi_n \begin{cases} \geq 0 & \text{if } t_n = 1 \\ < 0 & \text{otherwise} \end{cases}$$

- However, from  $\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n$  we see that the negative log likelihood is minimized when:  $y_n = \sigma(\mathbf{w}^T \phi_n) = t_n$ ,  $n = 1, \dots, N$ . This requires that the sigmoid function saturates which occurs when its argument

$$\mathbf{w}^T \phi_n \rightarrow \pm\infty, \text{ i.e. when } \|\mathbf{w}\| \rightarrow \infty.$$

# Logistic Regression- Linearly Separable Data

- MLE can exhibit severe over-fitting for data sets *that are linearly separable*.
- Indeed for linearly separable data, the hyperplane corresponding to  $P(C_1|x) = \sigma = 0.5 \rightarrow \mathbf{w}^T \phi = 0$  separates the two classes and will have the following:

$$\mathbf{w}^T \phi_n = \begin{cases} \geq 0 & \text{if } t_n = 1 \\ < 0 & \text{if } t_n = 0 \end{cases}$$

- In addition, we can see that the negative log-likelihood is minimized when  
 $\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n = 0 \Rightarrow y_n = \sigma(\mathbf{w}^T \phi_n) = t_n \quad \forall n, \text{i.e. } y_n = t_n = 0 \text{ or } y_n = t_n = 1$
- This occurs when *the sigmoid function is saturated*, which occurs when its argument  $\mathbf{w}^T \phi$  goes to  $\pm\infty$ , i.e. *when the magnitude of w goes to infinity*.
- The logistic sigmoid function becomes infinitely steep in feature space (step function) so every training point from each class  $k$  is assigned a posterior probability  $p(C_k | x) = 1$ .
- There are infinite solutions since any separating hyperplane will give the same posterior probabilities at the training data points.

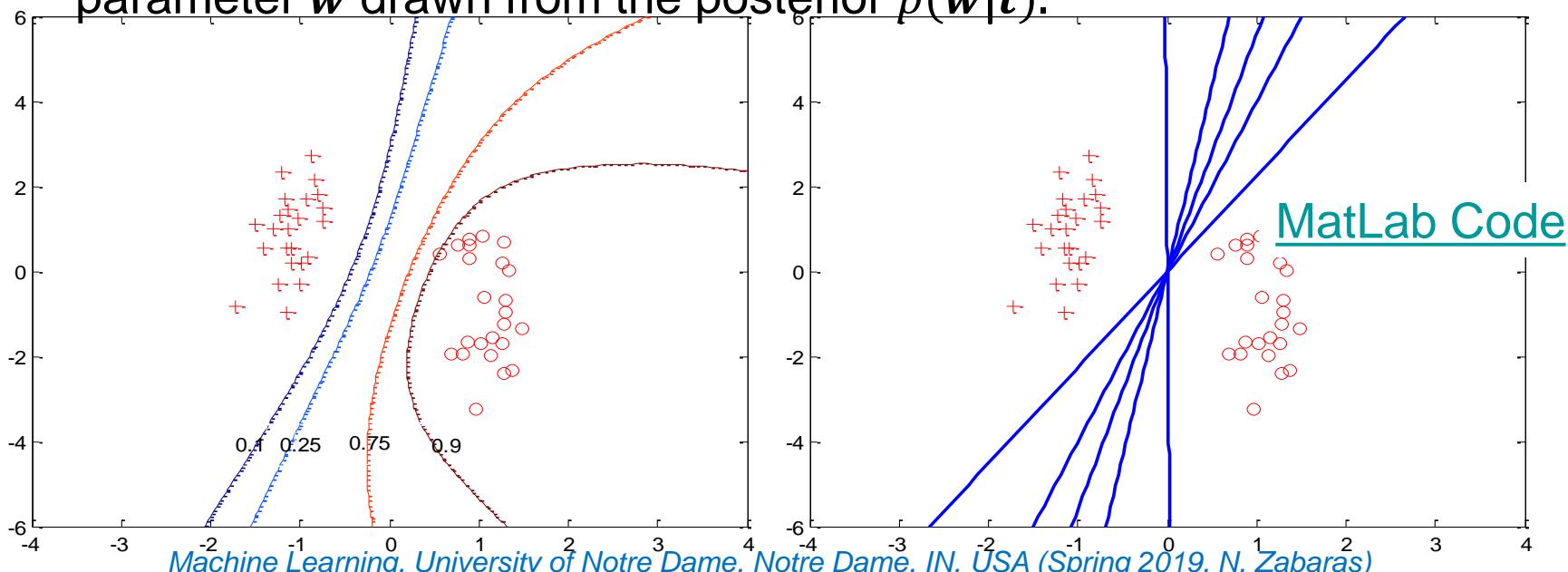
# **Logistic Regression Linearly Separable Data**

---

- Maximum likelihood provides no way to favor one such solution over another.
- *Which solution is found depends on the choice of optimization algorithm and on the parameter initialization.*
- The problem arises even if the number of data points is large compared to the number of parameters in the model, so long as the training data set is linearly separable.
- The singularity can be avoided by inclusion of a prior and finding a MAP solution for  $w$ , or equivalently by adding a regularization term to the error function.

# Logistic Regression-Linearly Separable Data

- Bayesian approach to logistic regression for linearly separable data set.
- Left: the predictive distribution using Variational Inference (details will be given in another lecture). The decision boundary lies between the clusters of data points. The contours of the predictive distribution splay out away from the data reflecting the greater uncertainty in the classification of such regions.
- Right: the decision boundaries corresponding to 5 samples of the parameter  $w$  drawn from the posterior  $p(w|t)$ .



# $L_2$ Regularization

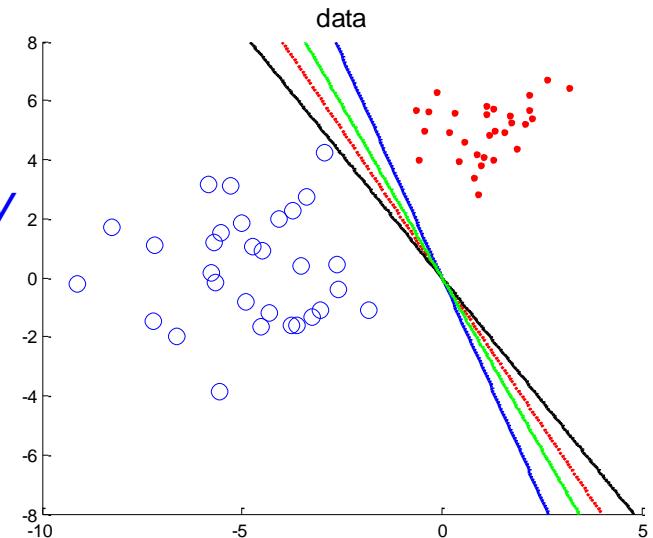
- ❑ Regularization is important for classification even with lots of data.
- ❑ For linearly separable data, the MLE is obtained by  $\|w\| \rightarrow \infty$ , corresponding to an infinitely steep sigmoid. In the two class linearly separable data, we show the Log-likelihood.
- ❑ The line is drawn from the origin in the direction of the MLE (which is at  $\infty$ ). The numbers correspond to the lines on the top Fig.
- ❑ This solution is very brittle and does not generalize well.
- ❑ To prevent this, we use  $L_2$  regularization.

$$f(\mathbf{w}) = NLL(\mathbf{w}) + \lambda \mathbf{w}^T \mathbf{w} \text{ (objective function)}$$

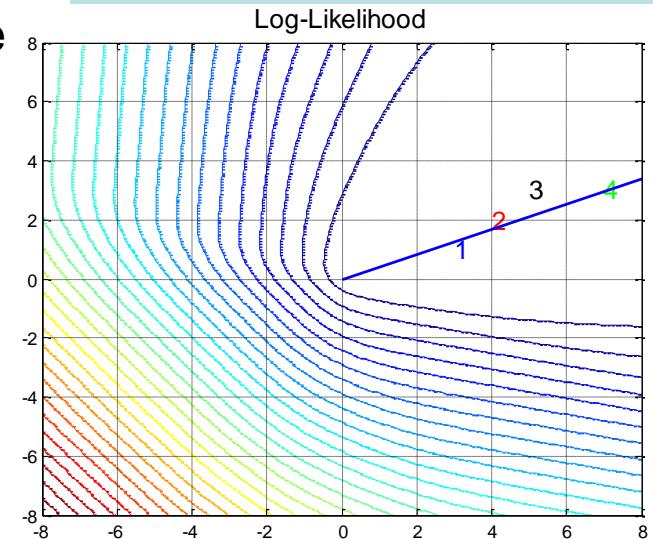
$$\mathbf{g}(\mathbf{w}) = \mathbf{g}(\mathbf{w}) + \lambda \mathbf{w} \text{ (gradient)}$$

$$\mathbf{H}(\mathbf{w}) = \mathbf{H}(\mathbf{w}) + \lambda \mathbf{I} \text{ (Hessian)}$$

- ❑ These modifications provide no difficulty with any gradient-based optimization toolbox.



[logregLaplaceGirolamiDemo](#)  
from [Kevin Murphys' PMTK](#)



# Iterative Reweighted Least Squares

- In linear regression models, the MLE for a Gaussian noise model had a closed-form solution. This was a consequence of the quadratic dependence of the log likelihood function on  $w$ .
- For logistic regression, there is no closed-form solution due to the nonlinearity of the logistic sigmoid function.
- However, the error function is convex and hence has a unique minimum.
- In addition, it can be minimized by an iterative technique based on *Newton-Raphson using a quadratic approximation to the log likelihood function*.
- The Newton-Raphson update, for minimizing a function  $E(w)$ , takes the form

$$w^{(new)} = w^{(old)} - H^{-1} \nabla E(w), \quad H_{ij} = \frac{\partial^2 E(w)}{\partial w_i \partial w_j} \text{ (Hessian)}$$

# IRLS for Linear Regression

- Let us first apply the Newton-Raphson method to the linear regression model with the sum-of-squares error function. The gradient and Hessian of this error function are given by

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^T \phi(x_n))^2 \Rightarrow \nabla E(\mathbf{w}) = \sum_{n=1}^N (\mathbf{w}^T \phi(x_n) - t_n) \phi_n = \Phi^T \Phi \mathbf{w} - \Phi^T \mathbf{t}$$

where  $\Phi$  is the  $N \times M$  matrix whose  $n^{\text{th}}$  row is given by  $\phi_n^T$ .

- The Newton-Raphson update takes the form:

$$\mathbf{w}^{(new)} = \mathbf{w}^{(old)} - (\Phi^T \Phi)^{-1} (\Phi^T \Phi \mathbf{w}^{(old)} - \Phi^T \mathbf{t}) = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

where  $\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \sum_{n=1}^N \phi_n \phi_n^T = \Phi^T \Phi$

- This is the standard least-squares solution to *linear regression*.
- The error function is quadratic and hence the Newton-Raphson gives the exact solution in one step!*

# IRLS Applied to Logistic Regression

- We similarly apply the Newton-Raphson to the cross-entropy error function for the logistic regression model:

$$E(\mathbf{w}) = -\ln p(\mathbf{t} | \mathbf{w}) = -\sum_{n=1}^N \left\{ t_n \ln y_n + (1-t_n) \ln (1-y_n) \right\}, \quad y_n = \sigma(a_n), \quad a_n = \mathbf{w}^T \boldsymbol{\phi}_n \Rightarrow$$

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \boldsymbol{\phi}_n = \boldsymbol{\Phi}^T (\mathbf{y} - \mathbf{t})$$

$$\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \sum_{n=1}^N y_n (1-y_n) \boldsymbol{\phi}_n \boldsymbol{\phi}_n^T = \boldsymbol{\Phi}^T \mathbf{R} \boldsymbol{\Phi}, \quad \mathbf{R} = \text{diag}(y_1(1-y_1), \dots, y_n(1-y_n))$$

where in the last equation we used  $d\sigma/da = \sigma(1-\sigma)$

- The Hessian now depends on  $\mathbf{w}$  through *the weighting matrix  $\mathbf{R}$* , since the error function is no longer quadratic.
- Using  $0 < y_n < 1$  from the form of the logistic sigmoid function, we see that  $\mathbf{u}^T \mathbf{H} \mathbf{u} > 0$  for any  $\mathbf{u}$ , and so  $\mathbf{H}$  is positive definite.
- The error function is a convex function of  $\mathbf{w}$  and hence has a unique minimum.

# IRLS Applied To Logistic Regression

- The Newton-Raphson update then becomes:

$$\begin{aligned}\mathbf{w}^{(new)} &= \mathbf{w}^{(old)} - (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T (\mathbf{y} - \mathbf{t}) = (\Phi^T \mathbf{R} \Phi)^{-1} \{ \Phi^T \mathbf{R} \Phi \mathbf{w}^{(old)} - \Phi^T (\mathbf{y} - \mathbf{t}) \} \\ &= (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} z, \text{ where } z = \Phi \mathbf{w}^{(old)} - \mathbf{R}^{-1} (\mathbf{y} - \mathbf{t})\end{aligned}$$

$$\mathbf{w} = \mathbf{0}_D;$$

$$w_0 = \log \bar{t} / (1 - \bar{t}) ;$$

repeat

$$\eta_i = w_0 + \mathbf{w}^T \mathbf{x}_i ;$$

$$y_i = \text{sigm}(\eta_i);$$

$$r_i = y_i(1 - y_i);$$

$$z_i = \eta_i + (t_i - y_i) / r_i;$$

$$\mathbf{R} = \text{diag}(r_{1:N});$$

$$\mathbf{w} = (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} z;$$

until converged;

- The update formula takes the form of a set of normal equations for a weighted least-squares problem.
- Because  $\mathbf{R}$  depends on  $\mathbf{w}$ , we apply the normal equations iteratively.
- For this reason, the algorithm is known as iterative reweighted least squares - IRLS.

- Rubin, D. B. (1983). [Iteratively reweighted least squares](#). In *Encyclopedia of Statistical Sciences*, Volume 4, pp. 272–275. Wiley.

# IRLS Applied To Logistic Regression

$$\mathbf{w}^{(new)} = (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} \mathbf{z}, \text{ where } \mathbf{z} = \Phi \mathbf{w}^{(old)} - \mathbf{R}^{-1} (\mathbf{y} - \mathbf{t}),$$

$$\mathbf{R} = \text{diag} (y_1(1-y_1), \dots, y_n(1-y_n))$$

- The elements of the diagonal weighting matrix  $\mathbf{R}$  can be interpreted as variances.
- Indeed note that the mean and variance of  $t$  in the logistic regression model are given as:

$$\mathbb{E}[t] = 1 \times \sigma(\mathbf{x}) + 0 \times (1 - \sigma(\mathbf{x})) = y$$

$$\text{var}[t] = \mathbb{E}[t^2] - \mathbb{E}[t]^2 = \sigma(\mathbf{x}) - \sigma(\mathbf{x})^2 = y(1-y)$$

where we used  $t^2 = t$  for  $t \in \{0, 1\}$ .

# IRLS Applied To Logistic Regression

$$\mathbf{w}^{(new)} = (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} \mathbf{z}, \text{ where } \mathbf{z} = \Phi \mathbf{w}^{(old)} - \mathbf{R}^{-1} (\mathbf{y} - \mathbf{t}),$$

$$\mathbf{R} = \text{diag} (y_1(1-y_1), \dots, y_n(1-y_n))$$

- We can also interpret IRLS as the solution to a linearized problem in the space of the variable  $a = \mathbf{w}^T \boldsymbol{\phi}$ .
- The quantity  $z_n$ , which corresponds to the  $n^{\text{th}}$  element of  $\mathbf{z}$ , can be seen as **an effective target value** in this space obtained by making a local linear approximation to the logistic sigmoid function around  $\mathbf{w}^{(old)}$ .

$$y_n = \sigma(a_n) = \frac{1}{1 + e^{-a_n}}, a_n = \mathbf{w}^T \boldsymbol{\phi}_n, \frac{dy_n}{da_n} = y_n(1 - y_n) \Rightarrow$$

$$t_n \cong y_n + \frac{dy_n}{da_n} \Big|_{\mathbf{w}^{(old)}} (a_n(w) - a_n(w^{(old)}))$$

$$a_n(w) \simeq a_n(w^{(old)}) + \frac{da_n}{dy_n} \Big|_{\mathbf{w}^{(old)}} (t_n - y_n) = \boldsymbol{\phi}_n^T \mathbf{w}^{(old)} - \frac{(y_n - t_n)}{y_n(1 - y_n)} = z_n$$

# Multiclass Logistic Regression

---

- For  $K > 2$  classes, the posterior probabilities are given by a *softmax transformation of linear functions of the feature variables* as follows:

$$p(C_k | \phi) = y_k(\phi) = \frac{e^{a_k}}{\sum_j e^{a_j}}, \quad a_k = \ln(p(\phi | C_k) p(C_k)) = \mathbf{w}_k^T \phi$$

- Use maximum likelihood to determine  $\{\mathbf{w}_k\}$  for this model.
- We need to compute the derivatives of  $y_k$  with respect to all of the activations  $a_j$ . They are given as:

$$y_k(\phi) = \frac{e^{a_k}}{\sum_j e^{a_j}} \Rightarrow \frac{\partial y_k}{\partial a_j} = y_k(\delta_{kj} - y_j)$$

- Krishnapuram, B., L. Carin, M. Figueiredo, and A. Hartemink (2005). Learning sparse bayesian classifiers: multi-class formulation, fast algorithms, and generalization bounds. IEEE Transaction on Pattern Analysis and Machine Intelligence.

# Multiclass Logistic Regression

---

- Next we write down the likelihood function.
- We use the 1-of-K coding scheme. The target  $t_n$  for a feature  $\phi_n$  belonging to class  $C_k$  is a binary vector with all elements zero except for element  $k$ , which equals one.
- The likelihood function is then given by

$$p(\mathbf{T} | \mathbf{w}_1, \dots, \mathbf{w}_K) = \prod_{n=1}^N \prod_{k=1}^K p(C_k | \phi_n)^{t_{nk}} = \prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}},$$

$$y_{nk} = y_k(\phi_n), \mathbf{T} \text{ } N \times K \text{ matrix}, T_{nk} = t_{nk}$$

- The negative log likelihood is then

$$E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\ln p(\mathbf{T} | \mathbf{w}_1, \dots, \mathbf{w}_K) = \sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk}$$

- This is the cross-entropy error function for the multiclass classification.

# Multiclass Logistic Regression

$$E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\ln p(\mathbf{T} \mid \mathbf{w}_1, \dots, \mathbf{w}_K) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk}$$

- We now take the gradient of the error function with respect to one of the  $\mathbf{w}_j$ .
- Using  $y_{nk} = \frac{e^{a_{kn}}}{\sum_j e^{a_{jn}}}$ ,  $\log y_{nk} = a_{kn} - \log \sum_j e^{a_{jn}}$ ,  $a_{kn} = \mathbf{w}_k^T \boldsymbol{\phi}_n$ ,  $\frac{\partial y_{nk}}{\partial a_j} = y_{nk} (\delta_{kj} - y_{nj})$ , and  $\sum_k t_{nk} = 1$ , we obtain

$$\nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = \sum_{n=1}^N (y_{nj} - t_{nj}) \boldsymbol{\phi}_n$$

- This is again the form for the gradient derived earlier for the sum-of-squares error function with the linear regression model and the cross-entropy error for logistic regression.
- We *use a sequential algorithm*

$$\mathbf{w}_j^{(\tau+1)} = \mathbf{w}_j^{(\tau)} - \eta \nabla E_{nj} = \mathbf{w}^{(\tau)} - \eta (y_{nj} - t_{nj}) \boldsymbol{\phi}_n$$

# General Form: Derivative of log Likelihood

- We have seen that the derivative of the log likelihood for a linear regression model with respect to  $w$  for a data point  $n$  took the form of the ‘error’  $y_n - t_n$  times the feature vector  $\phi_n$ .
- Similarly, for the combination of logistic sigmoid activation function and cross-entropy error function

$$E(w) = -\ln p(t | w) = -\sum_{n=1}^N \left\{ t_n \ln y_n + (1-t_n) \ln (1-y_n) \right\},$$

and for the softmax activation function with the multiclass cross-entropy error function

$$E(w_1, \dots, w_K) = -\ln p(T | w_1, \dots, w_K) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk}$$

- These are examples of a general result ([see here for more discussion](#))
  - Conditional distribution of the target variable from the exponential family
  - Choosing the activation function as the canonical link function

# IRLS Algorithm for the Multiclass Problem

- To find a batch algorithm, we again appeal to the Newton-Raphson update to obtain the corresponding IRLS algorithm for the multiclass problem.
- This requires evaluation of the Hessian matrix that comprises blocks of size  $M \times M$  in which block  $j, k$  is given by

$$\nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = \sum_{n=1}^N (y_{nj} - t_{nj}) \boldsymbol{\phi}_n \Rightarrow \nabla_{\mathbf{w}_k} \nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = \sum_{n=1}^N y_{nk} (\delta_{kj} - y_{nj}) \boldsymbol{\phi}_n \boldsymbol{\phi}_n^T$$

where we used:

$$\nabla_{\mathbf{w}_k} y_{nj} = \frac{dy_{nj}}{da_{nk}} \nabla_{\mathbf{w}_k} a_{nk} = y_{nk} (\delta_{kj} - y_{nj}) \boldsymbol{\phi}_n$$

- As with the two-class problem, the Hessian matrix for the multiclass logistic regression model is positive definite and so the error function again has a unique minimum.

- Bishop, C. M. and I. T. Nabney (2008). *Pattern Recognition and Machine Learning: A Matlab Companion*. Springer. In preparation

# MAP Estimate

- One can define an appropriate prior

$$p(\mathbf{W}) = \prod_k \mathcal{N}(\mathbf{w}_k | \mathbf{0}, \mathbf{V}_0)$$

- Using this, the modified negative log-likelihood and its gradient and Hessian are:

$$f'(\mathbf{W}) = -\log p(\mathcal{D}|\mathbf{W}) - \log p(\mathbf{W}) = f(\mathbf{W}) + \frac{1}{2} \sum_k \mathbf{w}_k \mathbf{V}_0^{-1} \mathbf{w}_k$$

$$\mathbf{g}'(\mathbf{W}) = \mathbf{g}(\mathbf{W}) + \mathbf{V}_0^{-1} \sum_k \mathbf{w}_k$$

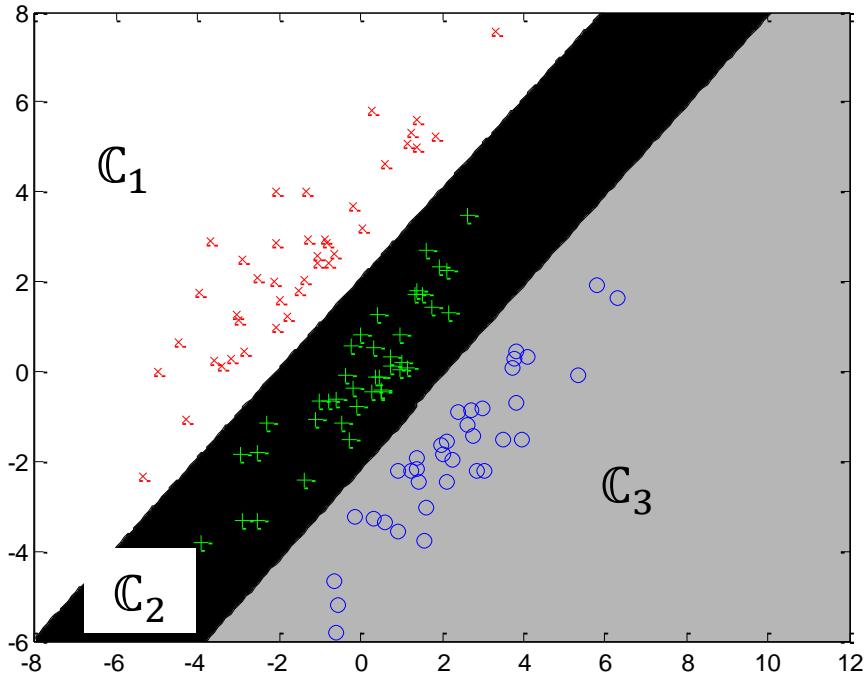
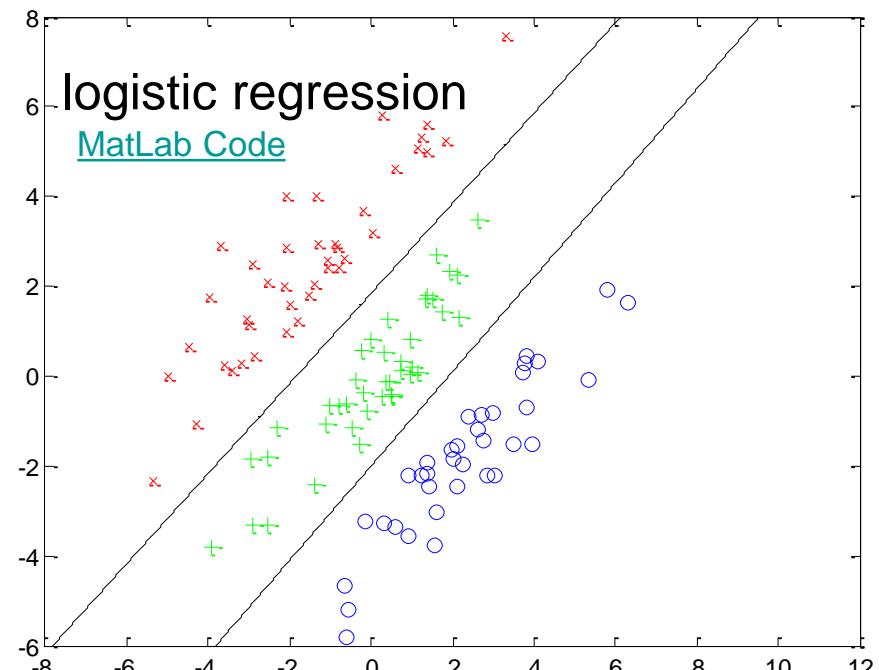
*logregFit*  
from [Kevin Murphys' PMTK](#)

$$\mathbf{H}'(\mathbf{W}) = \mathbf{H}(\mathbf{W}) + \mathbf{I}_K \otimes \mathbf{V}_0^{-1} ((KD) \times (KD))$$

- The symbol  $\mathbf{I}_C \otimes \mathbf{V}_0^{-1}$  is used here to denote [Kronecker matrix product](#).
- One can apply [\*limited memory BFGS\*](#) to find the MAP estimate.

# Multiclass Logistic Regression Example

This is the same example of 3 class classification discussed earlier.



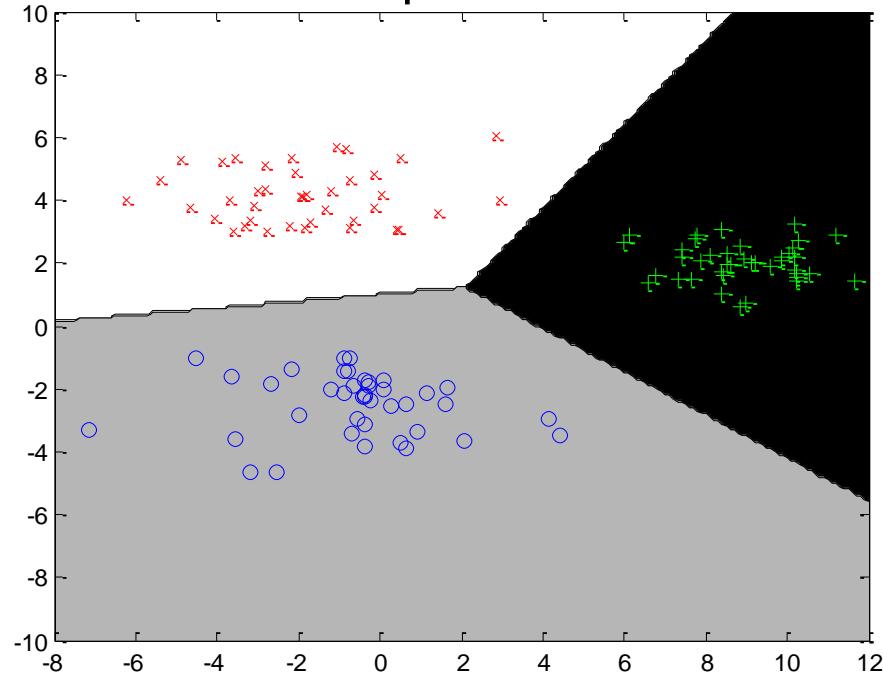
Here,  $a_1, a_2, a_3$  are the feature functions computed by following the logistic regression algorithm. Region  $C_1, C_2, C_3$  in the right plot can be interpreted as:

$$C_m: p(C_m|\phi) = \max(p(C_1|\phi), p(C_2|\phi), p(C_3|\phi)) \text{ where}$$

$$p(C_k|\phi) = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

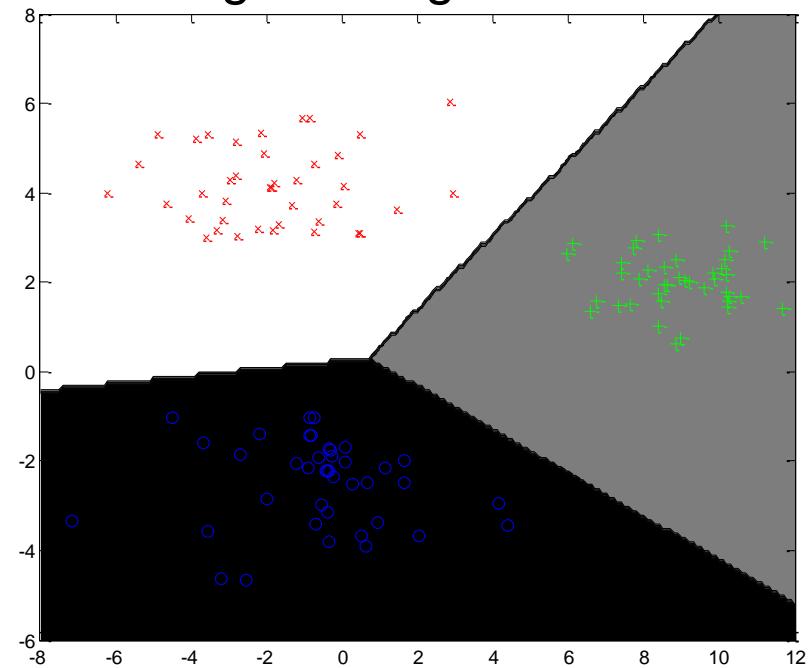
# Multiclass Logistic Regression Example

Least-squares



MatLab Code:  
Use 'data\_3c\_2'

Logistic Regression



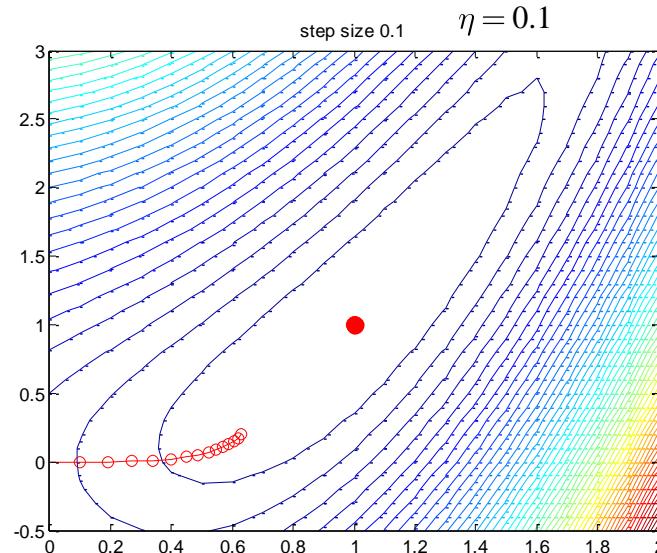
MatLab Code:  
Use 'data\_3c\_2'

---

# *Algorithms for Unconstrained Optimization*

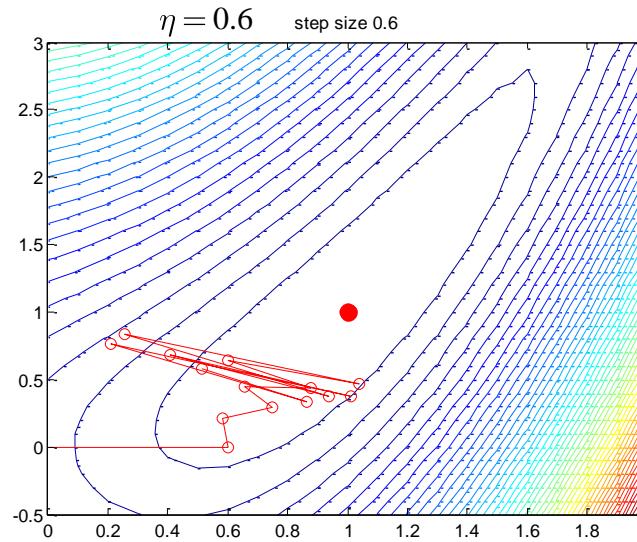
# Steepest Descent Algorithm

- The effect of the learning rate  $\eta$  on the performance of a steepest descent algorithm is shown on a simple function.



$$f(\theta) = 0.5 \theta_1^2 - \theta_2^2 + 0.5 \theta_1 - 1^2$$

$$\theta_{k+1} = \theta_k - \eta_k g_k = \theta_k - \eta_k \nabla f(\theta)|_k$$

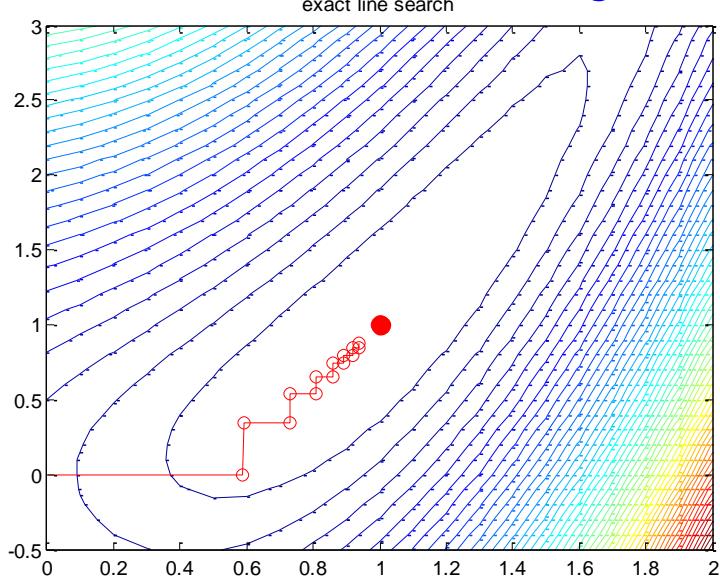


[steepestDescentDemo](#)  
from [Kevin Murphys' PMTK](#)

- Small step size leads to slow convergence and large step size to never converging. In both cases, the step size is kept constant.

# Gradient Descent Algorithm

- Gradient descent algorithm with line search is shown below.



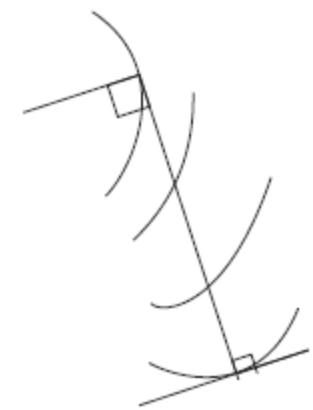
$$f(\boldsymbol{\theta}) = 0.5 \theta_1^2 - \theta_2^2 + 0.5 \theta_1 - 1^2$$

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \eta_k \mathbf{g}_k, \mathbf{g} = \nabla f$$

Pick  $\eta$  to minimize :

$$\phi(\eta) = f(\boldsymbol{\theta}_k) + \eta \mathbf{d}_k$$

$$\approx f(\boldsymbol{\theta}_k) + \eta \nabla f(\boldsymbol{\theta}_k) + \eta \mathbf{d}_k^T \mathbf{d}_k$$

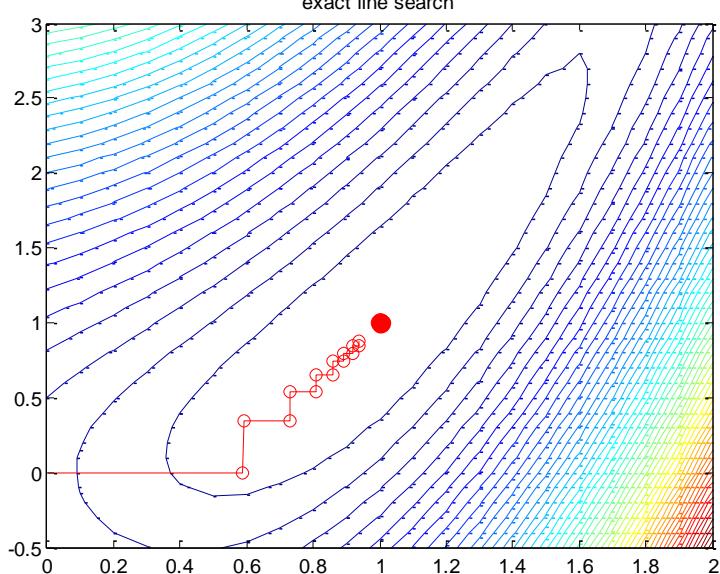


[steepestDescentDemo](#)  
from [Kevin Murphys' PMTK](#)

- Consecutive directions are orthogonal: noticing that  $\phi'(\eta) = 0$  leads to:  $\mathbf{d}^T \mathbf{g}_{k+1} = 0, \mathbf{g}_{k+1} = \nabla f(\boldsymbol{\theta}_k) + \eta \mathbf{d}_k$ .
- So at the end of each step  $k$ , either  $\nabla f|_{\boldsymbol{\theta}_{k+1}} = 0$  (stationary point was achieved) or  $\mathbf{d}_k \perp \nabla f|_{\boldsymbol{\theta}_{k+1}}$  (i.e. the search stops at a point where the gradient is normal to the search direction).

# Gradient Descent Example

- Gradient descent algorithm with line search is shown below.



[steepestDescentDemo](#)  
from [Kevin Murphys' PMTK](#)

- To avoid the zig-zag path, modify as:

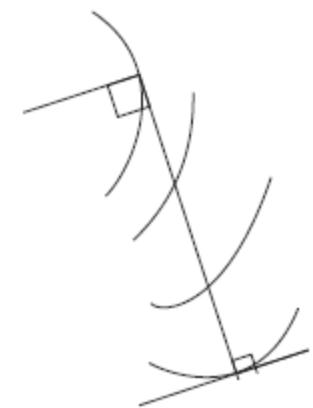
$$f(\theta) = 0.5 \theta_1^2 - \theta_2^2 + 0.5 \theta_1 - 1^2$$

$$\theta_{k+1} = \theta_k - \eta_k g_k, g = \nabla f$$

Pick  $\eta$  to minimize :

$$\phi(\eta) = f(\theta_k + \eta d_k)$$

$$\approx f(\theta_k) + \eta \nabla f(\theta_k)^\top d_k$$



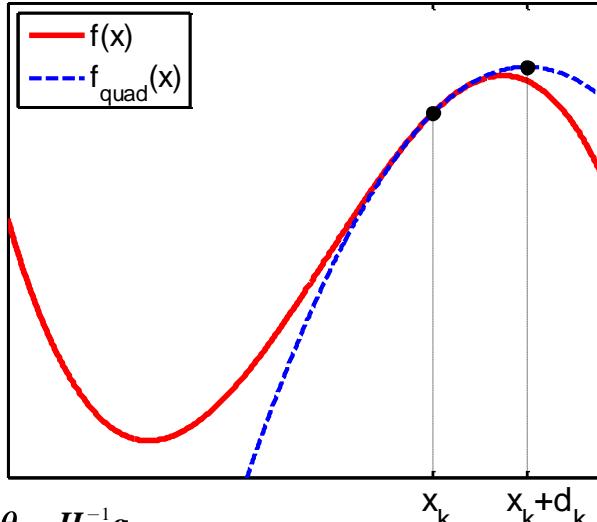
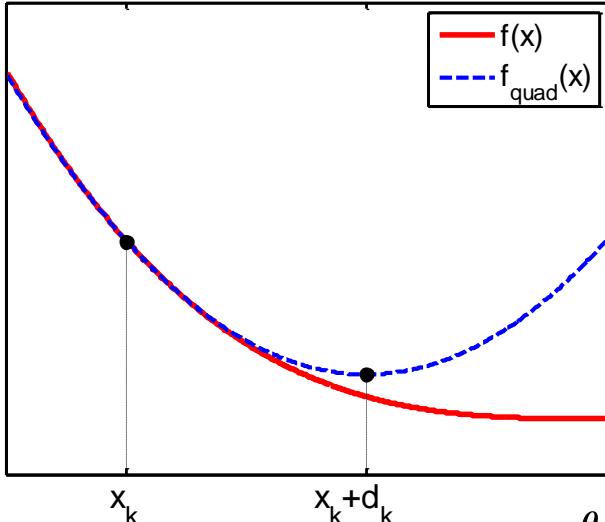
- Bertsekas, D. (1999). [Nonlinear Programming](#) (Second ed.). Athena Scientific.
- Nocedal, J. and S. Wright (2006). [Numerical Optimization](#). Springer.
- Golub, G. and C. F. van Loan (1996). [Matrix computations](#). Johns Hopkins University Press.

$$\theta_{k+1} = \theta_k - \eta_k g_k + \mu_k (\theta_k - \theta_{k-1}), \mu_k \in [0, 1]$$

- This is the so called *heavy-ball method*.
- Alternatively, one can apply the *conjugate gradient method*.

# Newton's Method

- Newton's method takes the Hessian into account:  $\theta_{k+1} = \theta_k - \eta_k \mathbf{H}_k^{-1} \mathbf{g}_k$ ,  $\mathbf{g} = \nabla f$



[newtonsMethodMinQuad](#)  
and [newtonsMethodNonConvex](#)  
from [Kevin Murphys' PMTK](#)

$$\theta_{k+1} = \theta_k \underbrace{- \mathbf{H}_k^{-1} \mathbf{g}_k}_{d_k}$$

- The step size  $d_k = -\mathbf{H}_k^{-1} \mathbf{g}_k$  is what should be added to  $\theta_k$  to minimize the 2<sup>nd</sup> order approximation of  $f$  around  $\theta_k$ .

$$\begin{aligned} f_{\text{quad}}(\theta) &= f_k + \mathbf{g}_k^T (\theta - \theta_k) + \frac{1}{2} (\theta - \theta_k)^T \mathbf{H}_k (\theta - \theta_k) \\ &= \theta^T A \theta + \mathbf{b}^T \theta + c \end{aligned}$$

$$A = \frac{1}{2} \mathbf{H}_k, \mathbf{b} = \mathbf{g}_k - \mathbf{H}_k \theta_k,$$

$$c = f_k - \mathbf{g}_k^T \theta_k + \frac{1}{2} \theta_k^T \mathbf{H}_k \theta_k$$

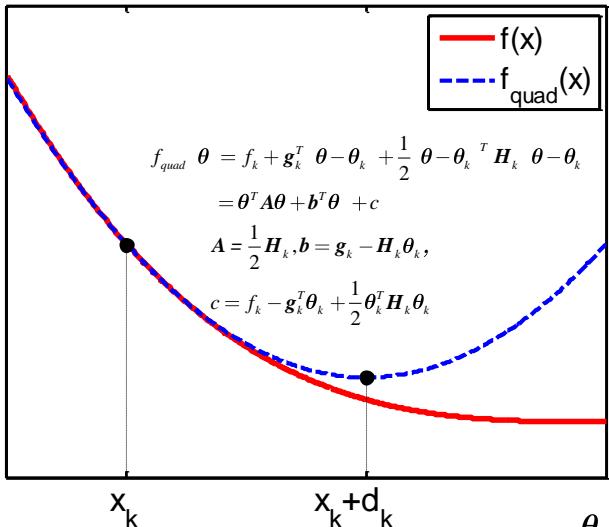
$$f_{\text{quad}}(\theta) = 0 \Rightarrow A\theta = -\frac{1}{2}\mathbf{b}$$

$$\Rightarrow \mathbf{H}_k \theta_{k+1} = -\mathbf{g}_k - \mathbf{H}_k \theta_k$$

$$\Rightarrow \theta_{k+1} = \theta_k - \mathbf{H}_k^{-1} \mathbf{g}_k$$

# Newton's Method

- Newton's method takes the Hessian into account:  $\theta_{k+1} = \theta_k - \eta_k \mathbf{H}_k^{-1} \mathbf{g}_k$ ,  $\mathbf{g} = \nabla f$



Initialize  $\theta_0$ ;

for  $k = 1, 2, \dots$  until convergence do

    Evaluate  $\mathbf{g}_k = \nabla f(\theta_k)$ ;

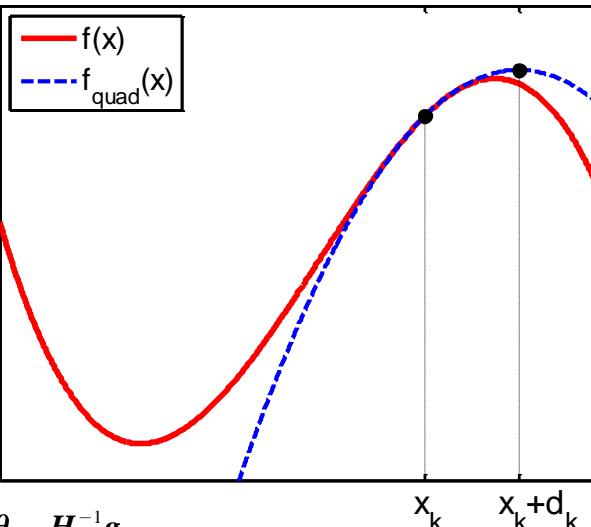
    Evaluate  $\mathbf{H}_k = \nabla^2 f(\theta_k)$ ;

    Solve  $\mathbf{H}_k \mathbf{d}_k = -\mathbf{g}_k$  for  $\mathbf{d}_k$ ;

    Use line search to find stepsize  $\eta_k$  along  $\mathbf{d}_k$ ;

$\theta_{k+1} = \theta_k + \eta_k \mathbf{d}_k$ ;

- Vandenberge, L. (2006). [Applied numerical computing](#): Lecture notes.



[newtonsMethodMinQuad](#)  
 and [newtonsMethodNonConvex](#)  
 from [Kevin Murphys' PMTK](#)

- **$H$  needs to be pos. definite** (holds for convex functions).

- If not,  $\mathbf{d}_k$  may not be a descent direction (Fig. on right).
- For such cases, use the **Levenberg Marquardt method** (adaptively switch between Newton and Steepest descent updates  $\mathbf{d}_k = -\mathbf{g}_k$ )

- Also one can use the **truncated Newton method** (solve  $\mathbf{H}_k \mathbf{d}_k = -\mathbf{g}_k$  with CG but truncate the iterations when negative curvature is observed).

# Limited Memory BFGS

---

- Computing  $\mathbf{H}$  is expensive. Quasi Newton methods iteratively build approximations using the gradient vector at each step.
- The BFGS method builds an approximation  $\mathbf{B}$  (rank–two updates) to the Hessian as (we start with  $\mathbf{B}_0 = \mathbf{I}$ ):

$$\mathbf{B}_{k+1} = \mathbf{B}_k + \frac{\mathbf{y}_k \mathbf{y}_k^T}{\mathbf{y}_k^T \mathbf{s}_k} - \frac{\mathbf{B}_k \mathbf{s}_k \mathbf{B}_k \mathbf{s}_k^T}{\mathbf{s}_k^T \mathbf{B}_k \mathbf{s}_k}$$

$$\mathbf{s}_k = \boldsymbol{\theta}_k - \boldsymbol{\theta}_{k-1}$$

$$\mathbf{y}_k = \mathbf{g}_k - \mathbf{g}_{k-1}$$

- BFGS can be thought of as a diagonal plus low-rank approximation to the Hessian.

- Nocedal, J. and S. Wright (2006). *Numerical Optimization*. Springer.

# Limited Memory BFGS

---

- Often we provide updates to the inverse of the Hessian:

$$\mathbf{C}_{k+1} = \left( \mathbf{I} - \frac{\mathbf{s}_k \mathbf{y}_k^T}{\mathbf{y}_k^T \mathbf{s}_k} \right) \mathbf{C}_k \left( \mathbf{I} - \frac{\mathbf{s}_k \mathbf{y}_k^T}{\mathbf{y}_k^T \mathbf{s}_k} \right) + \frac{\mathbf{s}_k \mathbf{s}_k^T}{\mathbf{y}_k^T \mathbf{s}_k}, \mathbf{C}_k \approx \mathbf{H}_k^{-1}$$

$$\mathbf{s}_k = \boldsymbol{\theta}_k - \boldsymbol{\theta}_{k-1}$$

$$\mathbf{y}_k = \mathbf{g}_k - \mathbf{g}_{k-1}$$

- To avoid storing the Hessian or its inverse, *limited memory BFGS (L-BFGS)* is used where  $\mathbf{H}_k^{-1} \mathbf{g}_k$  is performed by a sequence of inner products with  $\mathbf{s}_k$  and  $\mathbf{y}_k$  using only the  $m$  most recent ( $\sim 20$ ) pairs of  $(\mathbf{s}_k, \mathbf{y}_k)$ .
- This requires only  $\mathcal{O}(mD)$  storage vs  $\mathcal{O}(D^2)$  for the full Hessian.

- Nocedal, J. and S. Wright (2006). *Numerical Optimization*. Springer.

---

# *Bayesian Logistic Regression*

# Laplace Approximation

---

- The Laplace approximation is used to find a Gaussian approximation to a probability density.
- We use the fact that the log of a Gaussian is a quadratic function of the  $M$ -dimensional variables  $z$ . Taking a Taylor expansion of  $\ln f(z)$  centered on the mode  $z_0$

$$\ln f(z) \approx \ln f(z_0) - \frac{1}{2}(z - z_0)^T A(z - z_0)$$

*A is the  $M \times M$  Hessian matrix (needs to be positive definite)*

$$A = -\nabla \nabla \ln f(z) \Big|_{z=z_0}$$

- Taking the exponential of this approximation:

$$f(z) \approx f(z_0) \exp \left\{ -\frac{1}{2}(z - z_0)^T A(z - z_0) \right\}$$

# Laplace Approximation

---

- We can easily normalize this approximation as follows:

$$q(z) = \frac{|A|^{1/2}}{(2\pi)^{M/2}} \exp \left\{ -\frac{1}{2} (z - z_0)^T A (z - z_0) \right\} = \mathcal{N}(z | z_0, A^{-1})$$

- To apply the above approximation,  $f(z)$  doesn't need to be normalized.
- $A$  is the precision matrix that needs to be positive definite, i.e.  $z_0$  needs to be a local maximum.
- For multimodal distributions, there are many Laplace approximations depending which mode you are taking the approximation about.
- In Bayesian estimation, the posterior becomes closer to a Gaussian for large number of data (as a result of CLT).

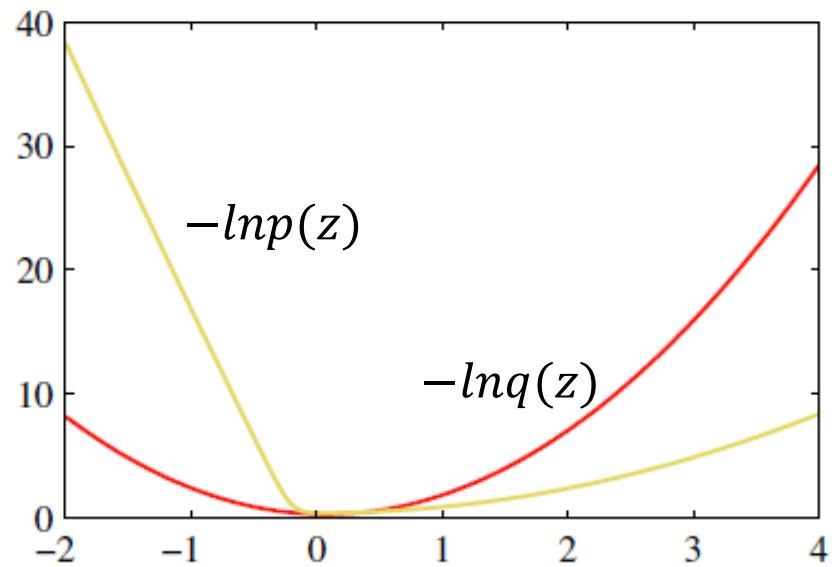
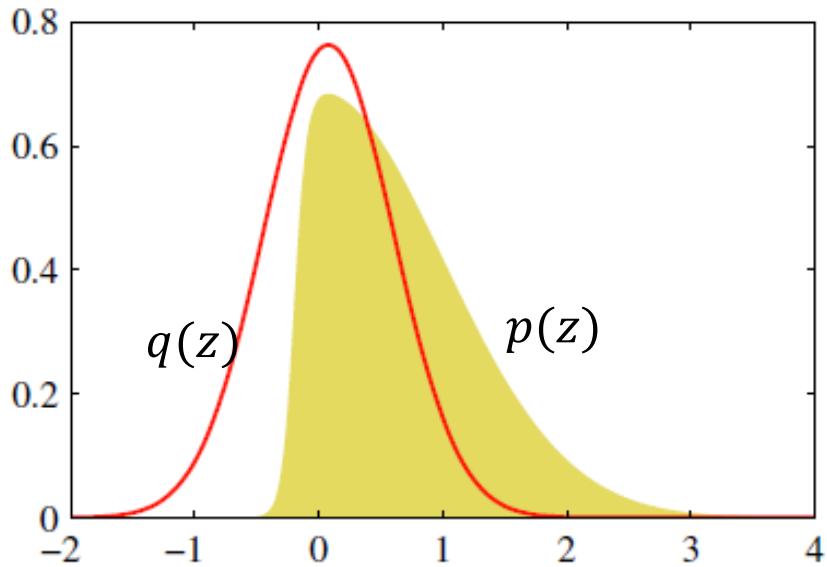
# **Laplace Approximation**

---

- Often we apply the Laplace approximation to a transformation of the variable. For instance if  $0 \leq \tau < \infty$  then we can consider a Laplace approximation of  $\ln\tau$ .
- Since the Laplace approximation is based on a Taylor expansion around the mode, it fails to capture important global properties.
- To capture global properties, alternative approximations are needed.

# Laplace Approximation

- The Laplace approximation applied to the distribution  $p(z) \propto \exp(-z^2/2)\sigma(20z + 4)$ ,  $\sigma(z)$  being the logistic sigmoid function  $\sigma(z) = (1 + e^{-z})^{-1}$ .
- The failure to capture global features is clear.



# Model Evidence and Model Comparison

---

- In addition to approximating  $p(\mathbf{z})$ , we can also *obtain an approximation to the normalization constant  $Z$* .

$$Z = \int f(\mathbf{z}) d\mathbf{z} = f(\mathbf{z}_0) \int \exp \left\{ -\frac{1}{2} (\mathbf{z} - \mathbf{z}_0)^T \mathbf{A} (\mathbf{z} - \mathbf{z}_0) \right\} d\mathbf{z} = f(\mathbf{z}_0) \frac{(2\pi)^{M/2}}{|\mathbf{A}|^{1/2}}$$

- We can *use this result to obtain an approximation to the model evidence needed in Bayesian model comparison.*
- Consider a data set  $\mathcal{D}$  and a set of models  $\{\mathcal{M}_i\}$  having parameters  $\{\boldsymbol{\theta}_i\}$ . For each model we define a likelihood function  $p(\mathcal{D} | \boldsymbol{\theta}_i, \mathcal{M}_i)$ .
- Introduce a prior  $p(\boldsymbol{\theta}_i | \mathcal{M}_i)$  over the parameters. We are interested in computing the model evidence  $p(\mathcal{D} | \mathcal{M}_i)$  for the various models. From now on we omit the conditioning on  $\mathcal{M}_i$  to keep the notation uncluttered.

# Model Evidence

- The model evidence is:

$$p(\mathcal{D}) = \int_Z \underbrace{p(\mathcal{D} | \boldsymbol{\theta}) p(\boldsymbol{\theta})}_{f(\boldsymbol{\theta})} d\boldsymbol{\theta}$$

- Using our earlier result, we can compute:

$$Z = \int f(z) dz = f(z_0) \frac{(2\pi)^{M/2}}{|A|^{1/2}} \Rightarrow$$
$$\ln p(\mathcal{D}) \simeq \underbrace{\ln p(\mathcal{D} | \boldsymbol{\theta}_{MAP})}_{Likelihood} + \underbrace{\ln p(\boldsymbol{\theta}_{MAP}) + \frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln |A|}_{\text{Occam factor}} \quad (\text{Penalizes Model Complexity})$$

- Here  $\boldsymbol{\theta}_{MAP}$  is the value of  $\boldsymbol{\theta}$  at the mode of the posterior distribution  $p(\boldsymbol{\theta} | \mathcal{D}) \propto p(\mathcal{D} | \boldsymbol{\theta}) p(\boldsymbol{\theta}) = f(\boldsymbol{\theta})$ , and  $A$  is the Hessian matrix of 2<sup>nd</sup> derivatives of the negative log posterior.

$$A = -\nabla \nabla \ln p(\mathcal{D} | \boldsymbol{\theta}_{MAP}) p(\boldsymbol{\theta}_{MAP}) = -\nabla \nabla \ln p(\boldsymbol{\theta}_{MAP} | \mathcal{D})$$

# Bayesian Information Criterion (BIC)

---

- Assuming that the Gaussian prior distribution over parameters is broad and that the Hessian has full rank, we can approximate the model evidence very roughly using

$$\begin{aligned}\ln p(\mathcal{D}) &\simeq \ln p(\mathcal{D}|\boldsymbol{\theta}_{MAP}) + \ln p(\boldsymbol{\theta}_{MAP}) + \frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{A}| \\ &\simeq \ln p(\mathcal{D}|\boldsymbol{\theta}_{MAP}) - \frac{1}{2} M \ln N\end{aligned}$$

- $N$  is the number of data points,  $M$  is the number of parameters in  $\boldsymbol{\theta}$  and we have omitted additive constants. This is the Bayesian Information Criterion (BIC). The proof is highlighted in the following slides.
- Compared to AIC (*Akaike information criterion*)

$$\ln p(\mathcal{D}|\widehat{\boldsymbol{\theta}}_{MLE}) - M$$

BIC penalizes model complexity more heavily.

- Akaike, H. (1974). [A new look at statistical model identification](#). *IEEE Trans. on Automatic Control* **19**, 716–723.
- Schwarz, G. (1978). [Estimating the dimension of a model](#). *Annals of Statistics* **6**, 461–464.

# Derivation of the BIC Criterion

---

- The BIC approximation can be viewed as a *large N approximation to the log model evidence*. We have:

$$\mathbf{A} = -\nabla \nabla \ln(p(|\boldsymbol{\theta}_{MAP}) p(\boldsymbol{\theta}_{MAP})) = \mathbf{H} - \nabla \nabla \ln p(\boldsymbol{\theta}_{MAP}), \quad \mathbf{H} = -\nabla \nabla \ln p(\mathcal{D} | \boldsymbol{\theta}_{MAP})$$

- If  $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta} | \mathbf{m}, \mathbf{V}_0)$ , the above equation is simplified:

$$\mathbf{A} = \mathbf{H} + \mathbf{V}_0^{-1}$$

from which we see that

$$\begin{aligned}\ln p(\mathcal{D}) &\simeq \ln p(\mathcal{D} | \boldsymbol{\theta}_{MAP}) + \ln p(\boldsymbol{\theta}_{MAP}) + \frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{A}| \Rightarrow \\ \ln p(\mathcal{D}) &\simeq \ln p(\mathcal{D} | \boldsymbol{\theta}_{MAP}) - \frac{1}{2} (\boldsymbol{\theta}_{MAP} - \mathbf{m})^T \mathbf{V}_0^{-1} (\boldsymbol{\theta}_{MAP} - \mathbf{m}) - \frac{1}{2} \ln |\mathbf{H}| + const\end{aligned}$$

# Derivation of the BIC Criterion

---

- Since we assume i.i.d. data,  $\mathbf{H} = -\nabla \nabla \ln p(\mathcal{D} | \boldsymbol{\theta}_{MAP})$  consists of a sum of terms, and we can consider the following approximation:

$$\mathbf{H} = \sum_{n=1}^N \mathbf{H}_n = N \widehat{\mathbf{H}}, \quad \widehat{\mathbf{H}} = \frac{\sum_{n=1}^N \mathbf{H}_n}{N}$$

- Combining this with the properties of the determinant, we have

$$\ln|\mathbf{H}| = \ln|N \widehat{\mathbf{H}}| = \ln(N^M |\widehat{\mathbf{H}}|) = M \ln N + \ln(|\widehat{\mathbf{H}}|)$$

where  $M$  is the dimensionality of  $\boldsymbol{\theta}$ . Note that we are assuming that  $\widehat{\mathbf{H}}$  has full rank  $M$ .

- We can now substitute the above equation in:

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D} | \boldsymbol{\theta}_{MAP}) - \frac{1}{2} (\boldsymbol{\theta}_{MAP} - \mathbf{m})^T \mathbf{V}_0^{-1} (\boldsymbol{\theta}_{MAP} - \mathbf{m}) - \frac{1}{2} \ln|\mathbf{H}| + const$$

# Derivation of the BIC Criterion

$$\ln p(\mathcal{D})$$

$$\simeq \ln p(\mathcal{D}|\boldsymbol{\theta}_{MAP}) - \frac{1}{2} (\boldsymbol{\theta}_{MAP} - \mathbf{m})^T V_0^{-1} (\boldsymbol{\theta}_{MAP} - \mathbf{m}) - \frac{1}{2} M \ln N - \frac{1}{2} \ln(|\hat{\mathcal{H}}|) + const$$

□ We finally obtain the required result by dropping the  $\ln(|\hat{\mathcal{H}}|)$

term since its  $\mathcal{O}(1)$  compared to  $\ln N$  and dropping the 2<sup>nd</sup> term with the assumption of prior with broad variance.

□ Thus finally we have:

$$BIC: \quad \ln p(\mathcal{D}) \simeq \ln p(\mathcal{D}|\boldsymbol{\theta}_{MAP}) - \frac{1}{2} M \ln N + const$$

# AIC and BIC

---

- Complexity measures such as AIC and BIC have the virtue of being easy to evaluate, but can also give misleading results.
- In particular, the assumption that the Hessian matrix has full rank is often not valid since many of the parameters are not ‘well-determined’.
- We can use

$$\ln p(\mathcal{D}) \simeq \underbrace{\ln p(\mathcal{D}|\boldsymbol{\theta}_{MAP})}_{Likelihood} + \underbrace{\frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln|A|}_{\begin{array}{c} Occam\ Factor \\ Penalizes\ Model\ Complexity \end{array}}$$

to obtain a more accurate estimate of the model evidence starting from the Laplace approximation.

# Bayesian Logistic Regression

- Exact Bayesian inference for logistic regression is intractable.
- Evaluation of the posterior distribution would require normalization of the product of a prior distribution and a likelihood function that is a product of logistic sigmoid functions one for every data point.
- Evaluation of the predictive distribution is similarly intractable.
- Here we consider the application of the Laplace approximation for the posterior. This requires evaluation of the 2<sup>nd</sup> derivatives of the log posterior (Hessian).

- [Spiegelhalter, D. and S. Lauritzen \(1990\). Sequential updating of conditional probabilities on directed graphical structures. Networks 20, 579–605.](#)
- [MacKay, D. J. C. \(1992b\). The evidence framework applied to classification networks. Neural Computation 4\(5\), 720–736.](#)
- Kuss and C. Rasmussen (2005). [Assessing approximate inference for binary gaussian process classification](#). *J. of Machine Learning Research* 6, 1679–1704.

# Bayesian Logistic Regression

- ❑ Because we seek a Gaussian representation for the posterior, we consider a Gaussian prior:

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0)$$

where  $\mathbf{m}_0, \mathbf{S}_0$  are fixed hyperparameters.

- ❑ The posterior distribution over  $\mathbf{w}$  is given as:

$$p(\mathbf{w} | \mathbf{t}) \propto p(\mathbf{w}) p(\mathbf{t} | \mathbf{w})$$

where  $\mathbf{t} = (t_1, \dots, t_N)^T$ . We take the log of both sides, and substitute for the prior distribution and for the likelihood function

$$p(\mathbf{t} | \mathbf{w}) = \prod_{n=1}^N y_n^{t_n} (1 - y_n)^{1-t_n}$$

,

# Bayesian Logistic Regression

$$\begin{aligned} \ln p(\mathbf{w} | \mathbf{t}) &\propto -\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1} (\mathbf{w} - \mathbf{m}_0) \\ &+ \sum_{n=1}^N \left\{ t_n \ln y_n + (1-t_n) \ln(1-y_n) \right\} + const., \end{aligned}$$

where  $y_n = \sigma(\mathbf{w}^T \boldsymbol{\phi}_n)$

- To obtain a Gaussian approximation to the posterior distribution, first maximize the posterior distribution to give the MAP solution  $\mathbf{w}_{MAP}$ , thus defining the mean of the Gaussian.
- The covariance is then given by the inverse of the matrix of 2<sup>nd</sup> derivatives of the negative log likelihood:

$$\mathbf{S}_N^{-1} = -\nabla \nabla \ln p(\mathbf{w} | \mathbf{t}) = \mathbf{S}_0^{-1} + \sum_{n=1}^N y_n (1-y_n) \boldsymbol{\phi}_n \boldsymbol{\phi}_n^T$$

# Gaussian Approximation to the Posterior

---

- The Gaussian approximation to the posterior now takes the form:

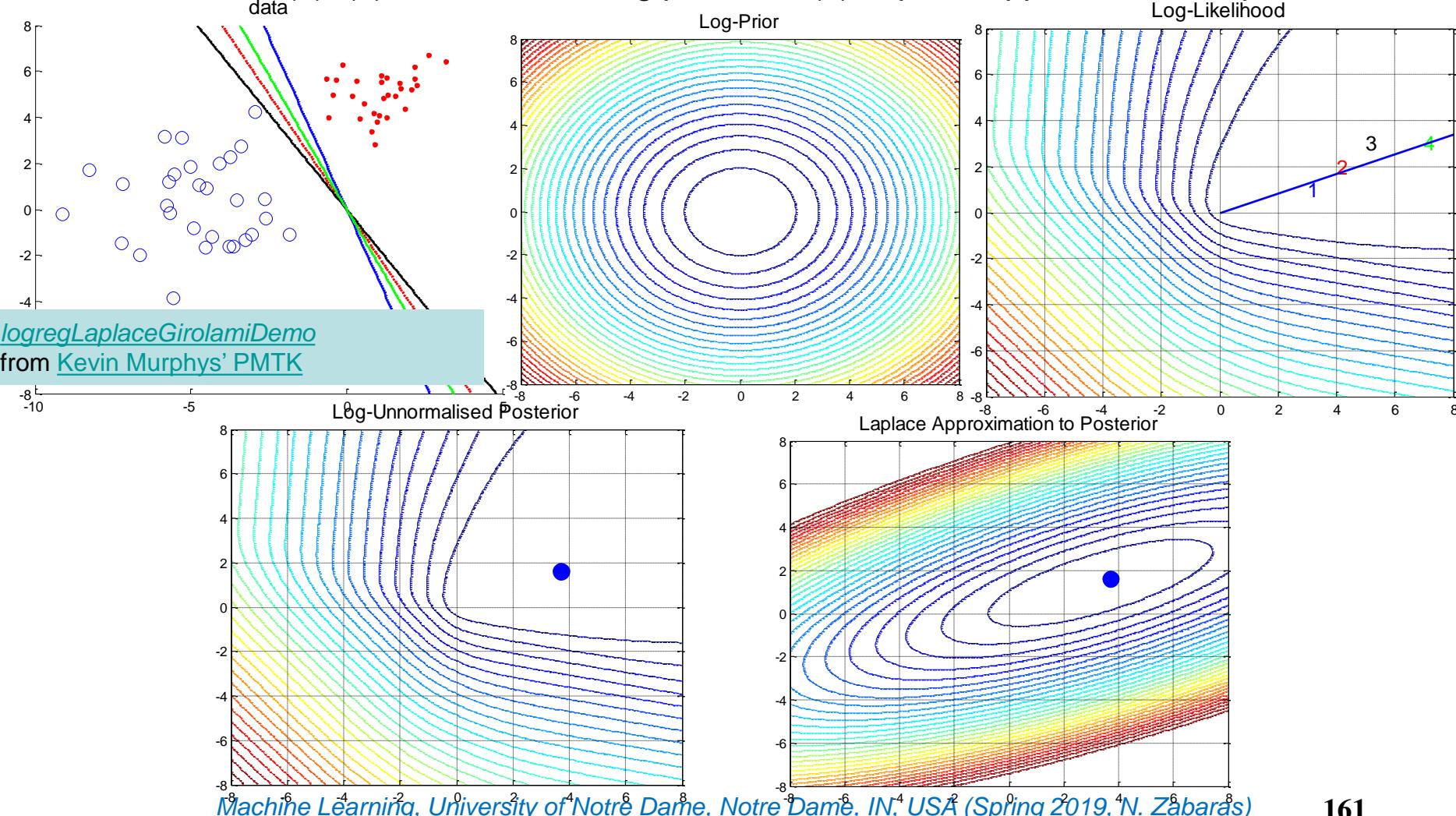
$$q(\boldsymbol{w}) = \mathcal{N}(\boldsymbol{w} \mid \boldsymbol{m}_{MAP}, \boldsymbol{S}_N)$$

$$\boldsymbol{S}_N^{-1} = -\nabla \nabla \ln p(\boldsymbol{w} \mid \boldsymbol{t}) = \boldsymbol{S}_0^{-1} + \sum_{n=1}^N y_n(1-y_n) \boldsymbol{\phi}_n \boldsymbol{\phi}_n^T$$

- To make predictions we need to be able to marginalize with respect to this distribution. That's not an easy task.

# Gaussian Approximation to the Posterior

- (a) Two-class data in 2d. (b) Log Prior  $\mathcal{N}(\mathbf{w}|0,100\mathbf{I})$  (c) Log-likelihood for a logistic regression model. The line shown is drawn from the origin in the direction of the MLE (which is at  $\infty$ ). The numbers correspond to the 4 points in  $\mathbf{w}$ -space corresponding to the lines in (a). (d) Unnormalized log posterior. (d) Laplace approximation to posterior.



# Predictive Distribution: Plug-In Approximation

---

- The posterior predictive distribution has the form

$$p(C_1 | \phi, t) = \int p(C_1 | \phi, w) p(w | t) dw$$

- The simplest approximation is the **plug-in approximation** (using the posterior mean) which for binary classification takes the form:

$$p(C_1 | \phi, t) = \int p(C_1 | \phi, w) p(w | t) dw \approx p(C_1 | \phi, \mathbb{E}[w])$$

- This approximation underestimates the uncertainty and better approximations are needed.

# Predictive Distribution: MC Approximation

---

- By taking samples from the posterior

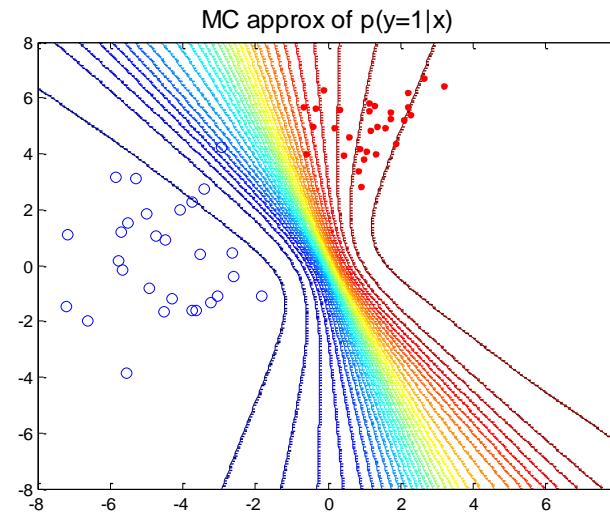
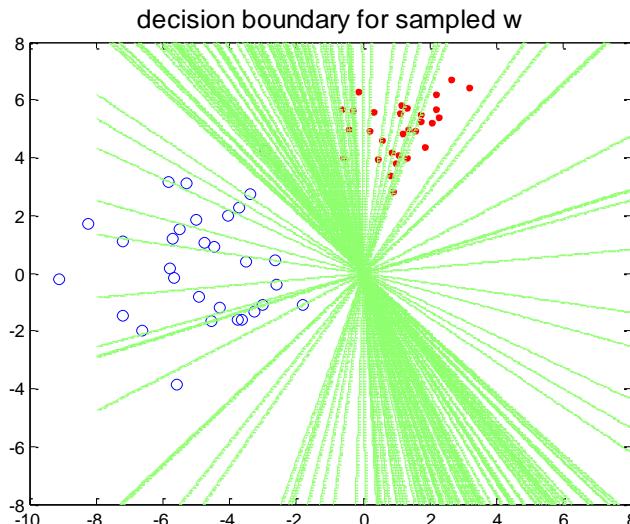
$$\mathbf{w}^s \sim p(\mathbf{w} | \mathbf{t})$$

$$p(C_1 | \phi, \mathbf{t}) = \int p(C_1 | \phi, \mathbf{w}) p(\mathbf{w} | \mathbf{t}) d\mathbf{w} \approx \frac{1}{S} \sum_{s=1}^S \sigma(\mathbf{w}^{sT} \phi(\mathbf{x}))$$

- This technique can be extended to the multi-class case.
- If we have approximated the posterior using MC, we can reuse these samples for prediction.
- If we made a Gaussian (Laplace) approximation to the posterior, we can draw *independent* samples from the Gaussian using standard methods.

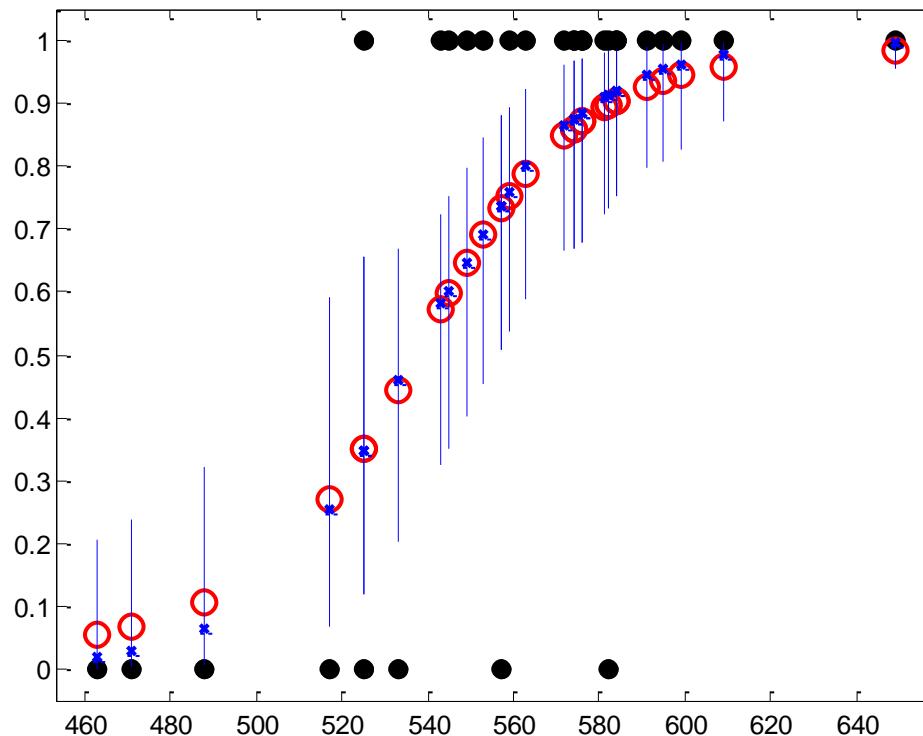
# Predictive Distribution: MC Approximation

- By averaging over multiple predictions, we see that the uncertainty in the decision boundary “splays out” as we move away from the training data.
- So although the decision boundary is linear, the posterior predictive density is not linear.
- Note also that the posterior mean decision boundary is roughly equally far from both classes.



# Predictive Distribution: MC Approximation

- *Red dots*: the mean of the posterior predictive evaluated at the training data.  
*Vertical blue lines*: 95% credible intervals for the posterior predictive; *Small blue star*: the median.
- With the Bayesian approach, we are able to model our uncertainty about the probability a student will pass the exam based on his SAT score.



[logregSATdemoBayes](#)  
from Kevin Murphys' PMTK

# Predictive Distribution: Probit Approximation

---

- The predictive distribution for class  $C_1$ , given a new feature vector  $\phi(x)$ , is obtained by marginalizing with respect to the *posterior distribution*  $p(w|t)$ , which is itself approximated by a Gaussian distribution  $q(w)$  so

$$p(C_1|\phi, t) = \int p(C_1|\phi, w) p(w|t) dw \simeq \int \sigma(w^T \phi) q(w) dw$$

The corresponding probability for class  $C_2$  is given simply by  
 $p(C_2|\phi, t) = 1 - p(C_1|\phi, t)$ .

- The function  $\sigma(w^T \phi)$  depends on  $w$  only through its projection onto  $\phi$ . Denoting  $a = w^T \phi$ , we have

$$\sigma(w^T \phi) = \int \delta(a - w^T \phi) \sigma(a) da \Rightarrow \int \sigma(w^T \phi) q(w) dw = \int \sigma(a) p(a) da$$
$$\int \delta(a - w^T \phi) q(w) dw$$

# Predictive Distribution: Probit Approximation

- We can evaluate  $p(a) = \int \delta(a - \mathbf{w}^T \boldsymbol{\phi}) q(\mathbf{w}) d\mathbf{w}$  by noting that the delta function imposes a linear constraint on  $\mathbf{w}$  and so forms a marginal distribution from the joint distribution  $q(\mathbf{w})$  by integrating out all directions orthogonal to  $\boldsymbol{\phi}$ .
- Because  $q(\mathbf{w})$  is Gaussian, the marginal distribution will also be Gaussian.
- We evaluate the mean and variance of this distribution by taking moments, and interchanging the order of integration over  $a$  and  $\mathbf{w}$ , so that

$$\mu_a = \mathbb{E}[a] = \int p(a) a da = \int q(\mathbf{w}) \mathbf{w}^T \boldsymbol{\phi} d\mathbf{w} = \mathbf{w}_{MAP}^T \boldsymbol{\phi}$$

$q(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{w}_{MAP}, \mathbf{S}_N)$

$$\begin{aligned}\sigma_a^2 &= \text{var}[a] = \int p(a) \left\{ a^2 - \mathbb{E}[a^2] \right\} da \\ &= \int q(\mathbf{w}) \left\{ (\mathbf{w}^T \boldsymbol{\phi})^2 - (\mathbf{w}_{MAP}^T \boldsymbol{\phi})^2 \right\} d\mathbf{w} = \boldsymbol{\phi}^T \mathbf{S}_N \boldsymbol{\phi}\end{aligned}$$

$q(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{w}_{MAP}, \mathbf{S}_N)$

# Predictive Distribution: Probit Approximation

---

- Note that the distribution  $p(a) = \mathcal{N}(a | \mathbf{w}_{MAP}^T \boldsymbol{\phi}, \boldsymbol{\phi}^T S_N \boldsymbol{\phi})$  takes the same form as the predictive distribution for the linear regression model with the noise variance set to zero.
- Thus our variational approximation to the predictive distribution is:

$$p(C_1 | \boldsymbol{\phi}, \mathbf{t}) = \int \sigma(a)p(a)da = \int \sigma(a)\mathcal{N}(a | \mu_a, \sigma_a^2)da$$
$$\Phi(\lambda a)$$

- We *approximate the logistic sigmoid function  $\sigma(a)$  with the inverse probit function  $\Phi(a)$ .*
- To obtain the best approximation to the logistic function we rescale the horizontal axis, so that we *approximate  $\sigma(a)$  by  $\Phi(\lambda a)$ . We find  $\lambda$  by requiring that the two functions have the same slope at the origin, which gives  $\lambda^2 = \pi/8$ .*

- [Spiegelhalter, D. and S. Lauritzen \(1990\). Sequential updating of conditional probabilities on directed graphical structures. Networks 20](#), 579–605.
- [MacKay, D. \(1992b\). The evidence framework applied to classification networks. Neural Comp. 4\(5\), 720–736.](#)
- Barber, D. and C. M. Bishop (1998a). [Ensemble learning for multi-layer networks](#). In M. I. Jordan, K. J. Kearns, and S. A. Solla (Eds.), [Advances in Neural Information Processing Systems](#), Volume 10, pp. 395–401.

# Predictive Distribution: Probit Approximation

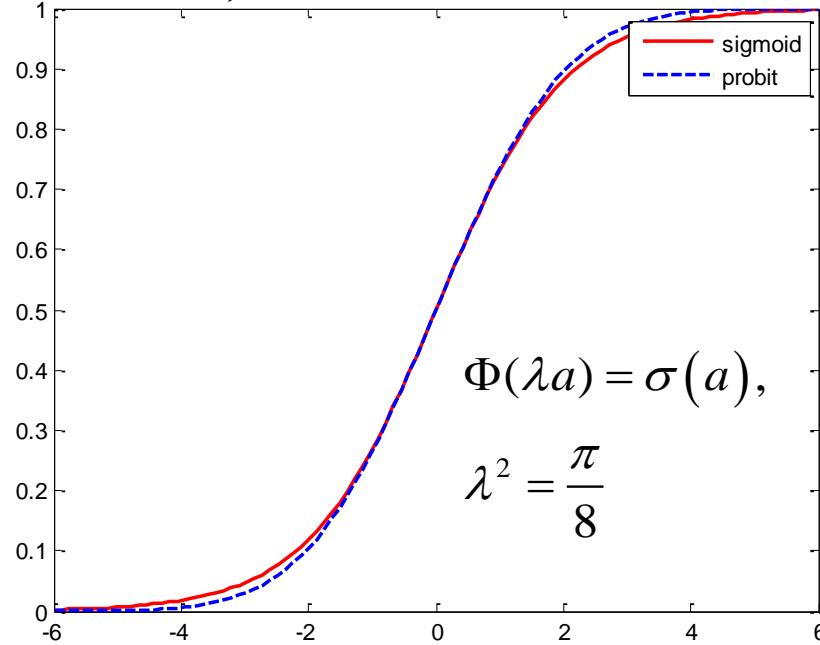
- The convolution below can be performed analytically:

$$\int \Phi(\lambda a) \mathcal{N}(a | \mu, \sigma^2) da = \Phi\left(\frac{\mu}{(\lambda^{-2} + \sigma^2)^{1/2}}\right)$$

where  $\Phi(a) = \int_{-\infty}^a \mathcal{N}(\theta | 0, 1) d\theta$

- We now apply the approximation  $\sigma(a) \simeq \Phi(\lambda a)$  to the inverse probit functions appearing on both sides of this equation, leading to the following approximation for the convolution of a logistic sigmoid with a Gaussian

$$\int \sigma(a) \mathcal{N}(a | \mu, \sigma^2) da = \sigma\left(\kappa(\sigma^2)\mu\right), \quad \kappa(\sigma^2) = \left(1 + \frac{\pi\sigma^2}{8}\right)^{-1/2}$$



*probitPlot*  
from [Kevin Murphys' PMTK](#)

# Appendix: Convolution of $\Phi(\lambda a)$ with a Gaussian

- We will highlight below the proof of

$$\int_{-\infty}^{\infty} \Phi(\lambda a) \mathcal{N}(a | \mu, \sigma^2) da = \Phi\left(\frac{\mu}{(\lambda^{-2} + \sigma^2)^{1/2}}\right), \text{ where } \Phi(b) = \int_{-\infty}^b \mathcal{N}(\theta | 0, 1) d\theta$$

- The derivative of the rhs wrt  $\mu$  is:  $\mathcal{N}\left(\frac{\mu}{(\lambda^{-2} + \sigma^2)^{1/2}} | 0, 1\right) \frac{1}{(\lambda^{-2} + \sigma^2)^{1/2}}$

- We use  $a = \mu + \sigma z$  and write the derivative wrt  $\mu$  on the lhs:

$$\begin{aligned} \int_{-\infty}^{\infty} \frac{d}{d\mu} \left\{ \Phi(\lambda(\mu + \sigma z)) \mathcal{N}((\mu + \sigma z) | \mu, \sigma^2) \right\} \sigma dz &= \int_{-\infty}^{\infty} \left\{ \frac{d\Phi(\lambda(\mu + \sigma z))}{d\mu} \mathcal{N}((\mu + \sigma z) | \mu, \sigma^2) + \Phi(\lambda(\mu + \sigma z)) \frac{d\mathcal{N}((\mu + \sigma z) | \mu, \sigma^2)}{d\mu} \right\} \sigma dz = \\ \int_{-\infty}^{\infty} \left\{ \mathcal{N}(\lambda(\mu + \sigma z) | 0, 1) \lambda \mathcal{N}((\mu + \sigma z) | \mu, \sigma^2) \right\} \sigma dz &= \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{-\frac{(\lambda(\mu + \sigma z))^2}{2}} e^{-\frac{z^2}{2}} \lambda dz \\ &\stackrel{M=\mu/\sqrt{\lambda^{-2}+\sigma^2}}{=} \frac{1}{\sqrt{2\pi}} e^{-M^2/2} \lambda \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{(\sqrt{\lambda^2\sigma^2+1}z+\lambda\sigma M)^2}{2}} dz = \\ \frac{1}{\sqrt{2\pi}} e^{-M^2/2} \lambda \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{\chi^2}{2}} \frac{1}{\sqrt{\lambda^2\sigma^2+1}} d\chi &= \frac{1}{\sqrt{2\pi}} \frac{1}{\sqrt{\sigma^2+\lambda^{-2}}} e^{-M^2/2} \end{aligned}$$

- The derivatives wrt  $\mu$  match. Thus the following must hold:

$$\int_{-\infty}^{\infty} \Phi(\lambda a) \mathcal{N}(a | \mu, \sigma^2) da = \Phi\left(\frac{\mu}{(\lambda^{-2} + \sigma^2)^{1/2}}\right) + \text{const}(\sigma, \lambda)$$

- Take  $\lambda = 0$  to see that the constant of integration is zero.

# Predictive Distribution: Probit Approximation

$$\int \sigma(a) \mathcal{N}(a | \mu, \sigma^2) da = \sigma\left(\kappa\left(\sigma^2\right)\mu\right), \quad \kappa\left(\sigma^2\right)=\left(1+\frac{\pi \sigma^2}{8}\right)^{-1/2}$$

□ Applying this result to

$$p(C_1 | t) = \int \sigma(a) p(a) da = \int_{\Phi(\lambda a)} \sigma(a) \mathcal{N}(a | \mu_a, \sigma_a^2) da$$

finally gives:

$$p(C_1 | \phi, t) = \sigma\left(\kappa\left(\sigma_a^2\right)\mu_a\right)$$

where  $\mu_a = \mathbf{w}_{MAP}^T \phi$  and  $\sigma_a^2 = \phi^T S_N \phi$

# Predictive Distribution: Probit Approximation

---

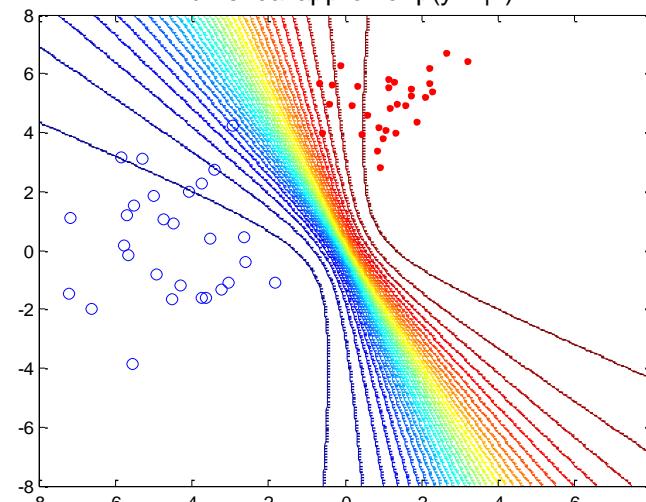
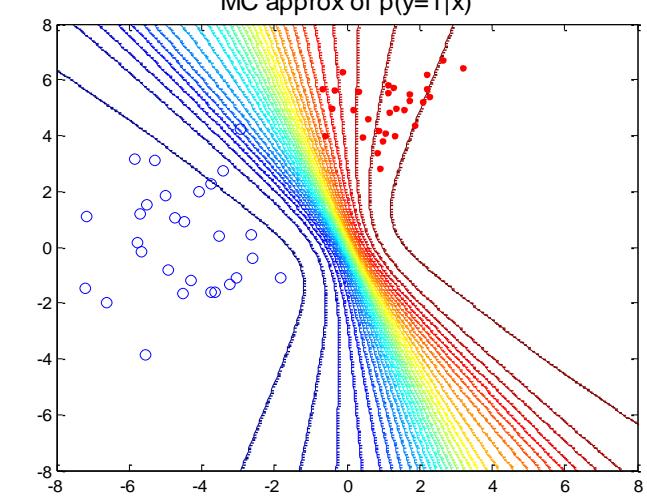
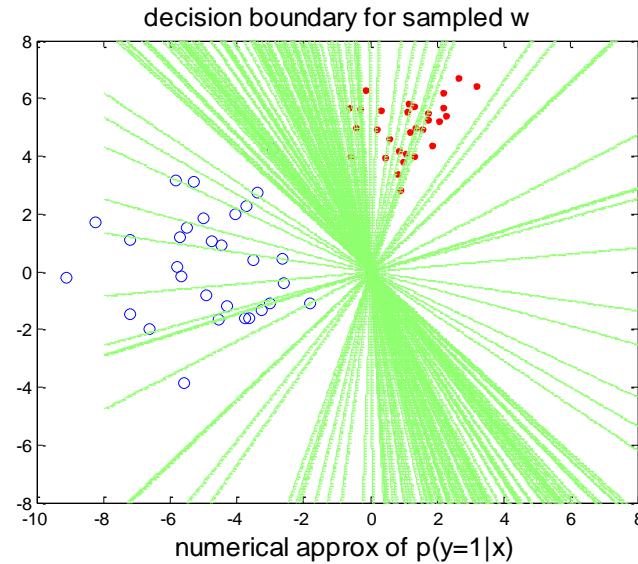
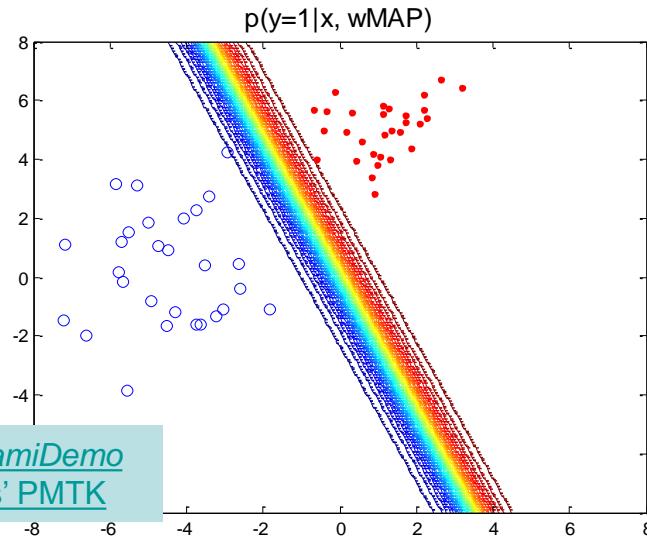
$$p(C_1 | \phi, t) = \sigma\left(\kappa\left(\sigma_a^2\right)\mu_a\right)$$

- Note that the decision boundary corresponding to  $p(C_1 | \phi, t) = 0.5$  is given by  $\mu_a = \mathbf{w}_{MAP}^T \phi = 0$ , which is the same as the decision boundary obtained by using the MAP value for  $\mathbf{w}$ .
- Thus for a decision criterion based on minimizing the misclassification rate and with equal prior probabilities, the marginalization over  $\mathbf{w}$  has no effect.
- For more complex decision criteria, using the predictive distribution is essential.

- Rasmussen, C. E. and C. K. I. Williams (2006). [Gaussian Processes for Machine Learning](#). MIT Press.

# Posterior Predictive

- Posterior predictive distribution for a logistic regression model in 2d. (a) contours of  $p(y = 1|x, \mathbf{w}_{map})$ . (b) samples from the posterior predictive distribution. (c) Averaging over these samples. (d) probit approximation.



# Outlier Detection

---

- It is useful to detect outliers in the data. In e.g. regression, this can be performed by computing

$$r_i = y_i - \hat{y}_i, \hat{y}_i = \hat{\mathbf{w}}^T \mathbf{x}_i$$

- These values should follow a  $\mathcal{N}(0, \sigma^2)$  distribution.
- This is assessed by a qq-plot, where we plot the  $N$  theoretical quantiles of a Gaussian distribution against the  $N$  empirical quantiles of the  $r_i$ .
- Points that deviate from the straight line are potential outliers.
- Methods based on residuals do not work for binary data since they rely on asymptotic statistics.

# *Outlier Detection in Logistic Regression*

---

- In logistic regression, we define outliers as points which have low probability.
- Outliers are defined as points with low probability under the cross-validated posterior predictive distribution,

$$p(t_i | \phi, \phi_{-i}, t_{-i}) = \int p(t_i | \phi_i, w) \underbrace{\prod_{i' \neq i} p(t_{i'} | \phi_{i'}, w)}_{\text{Posterior without point } i \text{ in likelihood}} p(w) dw$$

- This can be computed by sampling.

- Gelfand, A. (1996). [Model determination using sampling-based methods](#). In Gilks, Richardson, and Spiegelhalter (Eds.), *Markov Chain Monte Carlo in Practice*. Chapman & Hall.
- Johnson, V. and J. Albert (1999). [Ordinal data modeling](#). Springer.

---

# Probit Regression

# Probit Regression

---

- We have seen that for the exponential family of class-conditional distributions, the resulting posterior class probabilities are given by a logistic (or softmax) transformation acting on a linear function of the feature variables.
- Not all choices of class-conditional densities (e.g. Gaussian mixtures) give rise to this form for the posterior probabilities.
- We need to explore other types of discriminative probabilistic models.
- We return to the two-class case within the framework of generalized linear models so that

$$p(t = 1|a) = f(a)$$

where  $a = \mathbf{w}^T \boldsymbol{\phi}$ , and  $f(\cdot)$  is the activation function.

# Noisy Threshold Model

- We consider the generalized linear model:

$p(t = 1|a) = f(a)$ ,  $a = \mathbf{w}^T \boldsymbol{\phi}$ , and  $f(\cdot)$  is the activation function.

- To motivate an alternative choice for the link function  $f^{-1}(\cdot)$ , we consider a noisy threshold model, as follows. For each input  $\boldsymbol{\phi}_n$ , we evaluate  $a_n = \mathbf{w}^T \boldsymbol{\phi}_n$  and then we set the target value according to

$$\begin{cases} t_n = 1 \text{ if } a_n \geq \theta \\ t_n = 0 \text{ otherwise} \end{cases}$$

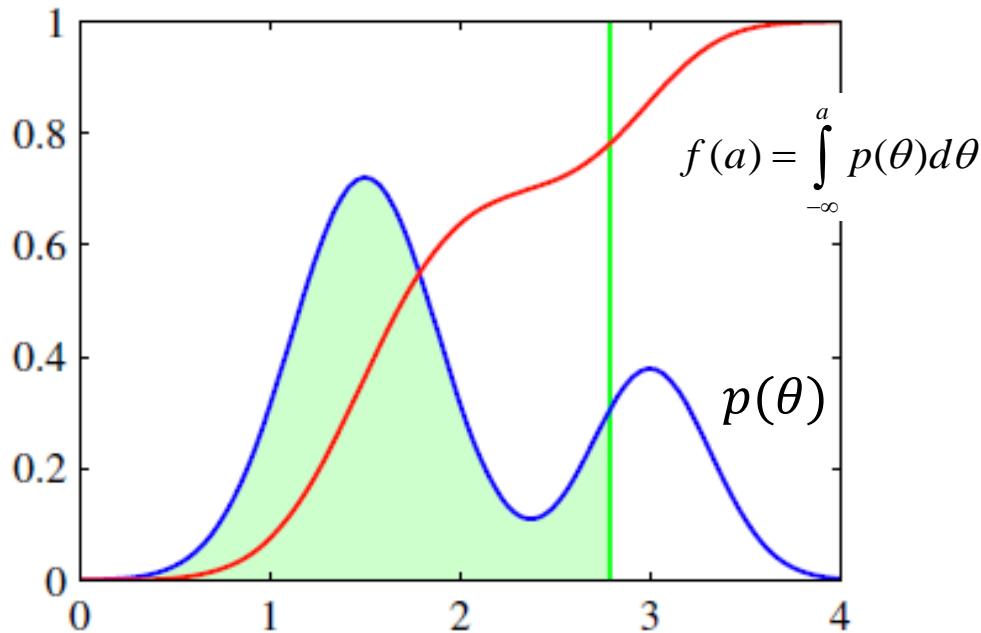
- The threshold  $\theta$  is stochastic and drawn from  $p(\theta)$ , the corresponding activation function will be given by the cumulative distribution function

$$f(a) = \int_{-\infty}^a p(\theta) d\theta$$

$$p(t=1|\boldsymbol{\phi}) \equiv p\left(\theta \leq \mathbf{w}^T \boldsymbol{\phi}\right) = \int_a^{\infty} p(\theta) d\theta = f(a)$$

# Noisy Threshold Model

- Here  $p(\theta)$  is a mixture of two Gaussians. Its cumulative distribution function  $f(a)$  is shown.
- In the stochastic threshold model, the class label takes the value  $t = 1$  if the value of  $a = \mathbf{w}^T \phi$  exceeds a threshold, otherwise it takes the value  $t = 0$ .
- This is equivalent to an activation function given by the cumulative distribution function  $f(a)$ .



# Inverse Probit Function

- As a specific example, suppose that the density  $p(\theta)$  is given by a zero mean, unit variance Gaussian. The corresponding cumulative distribution function is given by the inverse probit function:

$$\Phi(a) = \int_{-\infty}^a \mathcal{N}(\theta | 0, 1) d\theta$$

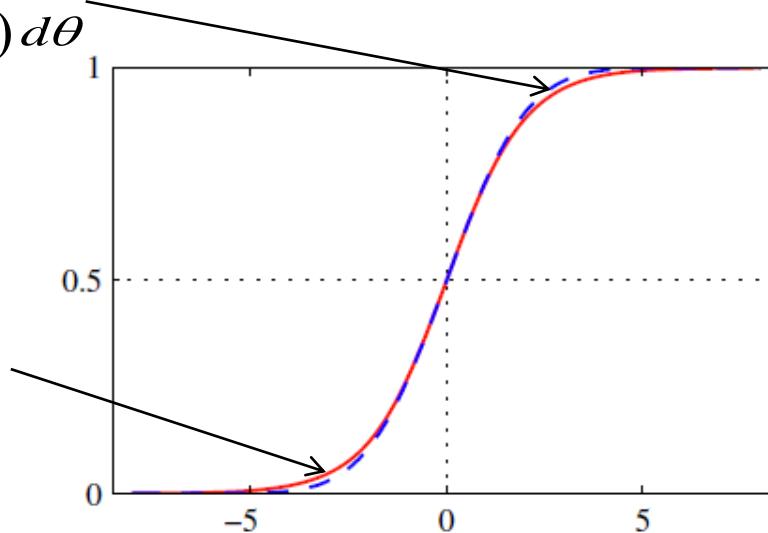
- It has a sinusoidal shape.

$$\Phi(\lambda a) = \int_{-\infty}^{\lambda a} \mathcal{N}(\theta | 0, 1) d\theta$$

Scaled inverse  
Probit function

$$\sigma(a) = \frac{1}{1 + e^{-a}}$$

Logistic  
sigmoid



[MatLab Code](#)

# Probit Regression

---

- ❑ Note that the use of a more general Gaussian distribution does not change the model because this is equivalent to a re-scaling of the linear coefficients  $w$ . Also note the relation:

$$erf(a) = \frac{2}{\sqrt{\pi}} \int_0^a e^{-\theta^2} d\theta \text{ (error function)}$$

$$\Phi(a) = \int_{-\infty}^a \mathcal{N}(\theta | 0, 1) d\theta \quad \Phi(a) = \frac{1}{2} \left\{ 1 + erf\left(\frac{a}{\sqrt{2}}\right) \right\}$$

- ❑ The generalized linear model based on a probit activation function is known as probit regression.
- ❑ We can determine the parameters of this model using MLE.
- ❑ The results using probit regression are similar to those of logistic regression.

# Outliers

---

- Outliers can arise through errors in measuring the input  $x$  or through mislabeling of the target value  $t$ .
- Such points can lie a long way to the wrong side of the decision boundary seriously distorting the classifier.
- The logistic & probit models behave differently to outliers:
  - the tails of the logistic sigmoid decays asymptotically  $\sim \exp(-x)$  for  $x \rightarrow \infty$ , whereas
  - for the probit activation function they decay  $\sim \exp(-x^2)$ .
- Thus the probit model is more sensitive to outliers.
- Mislabelling target values can easily be incorporated as follows:
$$p(t|\phi) = (1 - \epsilon) y(\phi) + \epsilon (1 - y(\phi)), \quad y(\phi) = \text{activation function}$$
- $\epsilon$  is the probability of mislabelling that can be set a priori or treated as hyperparameter to be inferred by the data.
  - [Opper, M. and O. Winther \(2000a\). Gaussian processes and SVM: mean field theory and leave-one-out](#). In A. J. Smola, P. L. Bartlett, B. Scholkopf, and D. Shuurmans (Eds.), *Advances in Large Margin Classifiers*, pp. 311–326. MIT Press.

# IRLS for Probit Regression

- To apply the IRLS iterative algorithm, we need to compute the gradient and Hessian of the log likelihood.
- The gradient is computed easily noting that:  $\nabla a_n = \phi_n$  and

$$E(\mathbf{w}) = -\ln p(\mathbf{t} | \mathbf{w}) = -\sum_{n=1}^N \left\{ t_n \ln y_n + (1-t_n) \ln(1-y_n) \right\} \Rightarrow \frac{\partial E}{\partial y_n} = \frac{y_n - t_n}{y_n(1-y_n)}$$

- Also from  $\Phi(a) = \frac{1}{2} \left\{ 1 + \operatorname{erf}\left(\frac{a}{\sqrt{2}}\right) \right\}$  and  $\operatorname{erf}(a) = \frac{2}{\sqrt{\pi}} \int_{-\infty}^a e^{-\theta^2} d\theta \Rightarrow$

$$\frac{\partial y_n}{\partial a_n} = \frac{\partial \Phi(a_n)}{\partial a_n} = \frac{1}{\sqrt{2\pi}} e^{-a_n^2/2}$$

- Thus the gradient is given as:

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N \frac{\partial E}{\partial y_n} \frac{\partial y_n}{\partial a_n} \nabla a_n = \sum_{n=1}^N \frac{y_n - t_n}{y_n(1-y_n)} \frac{1}{\sqrt{2\pi}} e^{-a_n^2/2} \phi_n$$

# IRLS for Probit Regression

- The computation of the Hessian is also straightforward. Note

$$\frac{\partial^2 E}{\partial y_n^2} = \frac{\partial}{\partial y_n} \frac{y_n - t_n}{y_n(1-y_n)} = \frac{y_n^2 + t_n^2 - 2y_n t_n}{y_n^2(1-y_n)^2}$$

- This leads to:

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N \frac{y_n - t_n}{y_n(1-y_n)} \frac{1}{\sqrt{2\pi}} e^{-a_n^2/2} \boldsymbol{\phi}_n \Rightarrow$$

$$\begin{aligned}\nabla^2 E(\mathbf{w}) &= \sum_{n=1}^N \left[ \frac{\partial}{\partial y_n} \left\{ \frac{y_n - t_n}{y_n(1-y_n)} \right\} \frac{1}{\sqrt{2\pi}} e^{-a_n^2/2} \boldsymbol{\phi}_n \nabla y_n + \frac{y_n - t_n}{y_n(1-y_n)} \frac{1}{\sqrt{2\pi}} e^{-a_n^2/2} (-a_n) \boldsymbol{\phi}_n \nabla a_n \right] \\ &= \sum_{n=1}^N \left[ \frac{y_n^2 + t_n^2 - 2y_n t_n}{y_n(1-y_n)} \frac{1}{\sqrt{2\pi}} e^{-a_n^2/2} - a_n (y_n - t_n) \right] e^{-a_n^2/2} \frac{1}{\sqrt{2\pi}} \frac{1}{y_n(1-y_n)} \boldsymbol{\phi}_n \boldsymbol{\phi}_n^T\end{aligned}$$

---

## Generalizing our Results:

(1)  $p(t|\phi(x)) = f(\mathbf{w}^T \phi(x))$  from the Exponential Family, with

(2)  $f$  being the Canonical Link Function

# Canonical Link Functions

- For the linear regression model with a Gaussian noise, the error function, corresponding to the negative log likelihood, is

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n))^2$$

- If we take the derivative wrt  $\mathbf{w}$  of the contribution to the error function from a data point  $n$ , this takes the form of the ‘error’  $y_n - t_n$  times the feature vector  $\boldsymbol{\phi}_n$ , where  $y_n = \mathbf{w}^T \boldsymbol{\phi}_n$ .
- Similarly, for the combination of the logistic sigmoid activation function and the cross-entropy error function

$$E(\mathbf{w}) = -\ln p(\mathbf{t} | \mathbf{w}) = -\sum_{n=1}^N \left\{ t_n \ln y_n + (1-t_n) \ln (1-y_n) \right\},$$

and for the softmax activation function with the multiclass cross-entropy error function

$$E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\ln p(\mathbf{T} | \mathbf{w}_1, \dots, \mathbf{w}_K) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk}$$

we again obtain this same simple form.

# Canonical Link Functions

---

- This is a general result: Assume a conditional distribution for the target variable from the exponential family, along with a corresponding choice for the activation function known as the canonical link function.
- We again make use of the restricted form of exponential family distributions.

$$p(\mathbf{x} | \boldsymbol{\lambda}_k, s) = \frac{1}{s} h\left(\frac{1}{s} \mathbf{x}\right) g(\boldsymbol{\lambda}_k) \exp\left\{\frac{1}{s} \boldsymbol{\lambda}_k^T \mathbf{x}\right\}$$

- Note that here we are applying the assumption of exponential family distribution to the target variable  $t$ , in contrast to earlier discussion where we applied it to the input vector  $\mathbf{x}$ .
- We thus consider conditional distributions of the target variable  $t$  of the form

$$p(t | \eta, s) = \frac{1}{s} h\left(\frac{t}{s}\right) g(\eta) \exp\left\{\frac{\eta t}{s}\right\}$$

# Canonical Link Functions

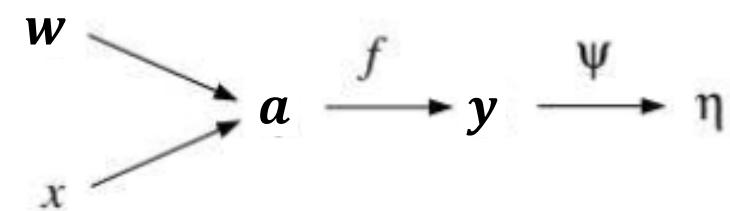
$$p(t | \eta, s) = \frac{1}{s} h\left(\frac{t}{s}\right) g(\eta) \exp\left\{\frac{\eta t}{s}\right\}$$

- We will prove in a [follow up lecture on generalized linear models](#) that:

$$-\nabla \ln g(\eta) = \mathbb{E}[u(x)] \text{ for } p(x | \eta) = h(x)g(\eta)\exp(\eta^T u(x))$$

we see that the conditional mean of  $t$ , which we denote by  $y$ , is

$$y \equiv \mathbb{E}[t | \eta] = -s \frac{d}{d\eta} \ln g(\eta)$$



- Thus  $y$  and  $\eta$  must be related:  $\eta = \psi(y)$
- We define a generalized linear model to be one for which  $y$  is a nonlinear function of a linear combination of the input (or feature) variables so that  $y = f(w^T \phi)$ .
- $f(\cdot)$  is the activation function, and  $f^{-1}(\cdot)$  is known as the link function in statistics.

- Nelder, J. A. and R.W. M. Wedderburn (1972). [Generalized linear models. Journal of the Royal Statistical Society, A 135, 370–384.](#)

# Canonical Link Functions

- Now consider the log likelihood for this model (as a function of  $\eta$ ):

$$\ln p(t | \eta, s) = \sum_{n=1}^N p(t_n | \eta, s) = \sum_{n=1}^N \left( \ln g(\eta_n) + \frac{\eta_n t_n}{s} \right) + const$$

where we are assuming that  $s$  is independent of  $\eta$ , i.e. all observations share a common scale parameter (which for a Gaussian distribution corresponds to the noise variance).

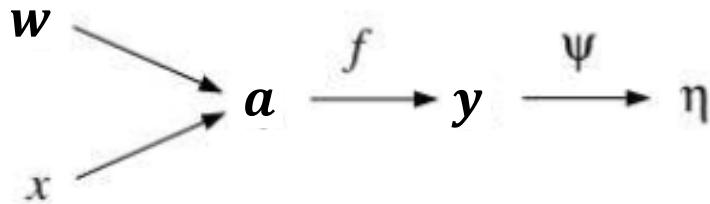
- The derivative of the log likelihood with respect to the model parameters  $w$  is then given by

$$\nabla_w \ln p(t | \eta, s) = \sum_{n=1}^N \left( \frac{d}{d\eta_n} \ln g(\eta_n) + \frac{t_n}{s} \right) \frac{d\eta_n}{dy_n} \frac{dy_n}{da_n} \nabla a_n = \frac{1}{s} \sum_{n=1}^N (t_n - y_n) \psi'(y_n) f'(a_n) \phi_n$$

where:  $a_n = w^T \phi_n$ ,  $y_n = f(a_n)$ ,  $y = \mathbb{E}[t | \eta] = -s \frac{d}{d\eta} \ln g(\eta) (\Rightarrow \eta = \psi(y))$

- There is considerable simplification if we take

$$f^{-1}(y) = \psi(y) \Rightarrow f(\psi(y)) = y \Rightarrow \psi'(y) f'(\psi) = 1$$



# Canonical Link Functions

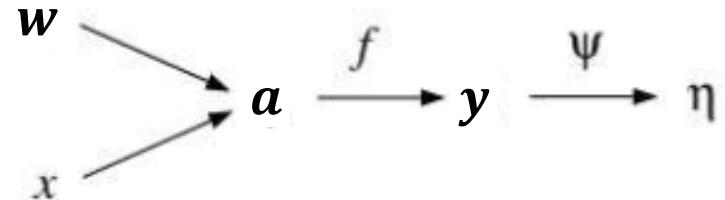
$$\nabla_w \ln p(t | \eta, s) = \frac{1}{s} \sum_{n=1}^N (t_n - y_n) \psi'(y_n) f'(a_n) \phi_n$$

where:  $a_n = \mathbf{w}^T \phi_n$ ,  $y_n = f(a_n)$ ,  $y = \mathbb{E}[t | \eta] = -s \frac{d}{d\eta} \ln g(\eta)$

- There is considerable simplification if we take

$$f^{-1}(y) = \psi(y) \Rightarrow f(\psi(y)) = y \Rightarrow \psi'(y)f'(\psi) = 1$$

- Also because



$$f^{-1}(y) = \psi(y) \Rightarrow a = \psi(y) \quad \text{Also use: } \psi'(y)f'(\psi) = 1 \Rightarrow \psi'(y)f'(a) = 1$$

- The gradient of the error is now simplified to a familiar form as:

$$\nabla E(\mathbf{w}) = \frac{1}{s} \sum_{n=1}^N (t_n - y_n) \phi_n, \quad \begin{array}{l} \text{Gaussian model: } s = \beta^{-1} \\ \text{Logistic model: } s = 1 \end{array}$$